

ΠΑΝΕΠΙΣΤΗΜΙΟ ΙΩΑΝΝΙΝΩΝ  
ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ Η/Υ &  
ΠΛΗΡΟΦΟΡΙΚΗΣ  
ΜΥΕ002 - ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ  
ΑΣΚΗΣΗ 2

Κίτσιος Κωνσταντίνος 4388  
Λιόντος Αναστάσιος 4409

June 20, 2021

## Contents

1	K-means clustering	3
2	Hierarchical clustering	7
3	K-medoids algorithm	12
4	Βοηθητική συνάρτηση histogram	16
5	Βοηθητική συνάρτηση reduce	16
6	Η βέλτιστη μέθοδος	16

# 1 K-means clustering

```
1 import idx2numpy as idx2numpy
2 import numpy as np
3 import Histogram
4 from sklearn.metrics import confusion_matrix
5
6
7 trainImagePath = "../dataset/train-images-idx3-ubyte"
8 trainLabelPath = "../dataset/train-labels-idx1-ubyte"
9
10 testImagePath = "../dataset/t10k-images-idx3-ubyte"
11 testLabelPath = "../dataset/t10k-labels-idx1-ubyte"
12
13 x_train = idx2numpy.convert_from_file(trainImagePath)
14 y_train = idx2numpy.convert_from_file(trainLabelPath)
15
16 x_test = idx2numpy.convert_from_file(testImagePath)
17 y_test = idx2numpy.convert_from_file(testLabelPath)
18
19 x_train=x_train.reshape((x_train.shape[0], 28*28))
20 x_test=x_test.reshape((x_test.shape[0], 28*28))
21
22 x_train=x_train/255.0
23 x_test=x_test/255.0
24
25 k=10; # number of clusters
26 centroids = {}
27 def fit(data, distance):
28
29     for i in range(k):
30         centroids[i] = data[i]
31
32     while(True):
33         classifications = {}
34
35         for i in range(k):
36             classifications[i] = []
37
38         for featureset in data:
39             if (distance=='Euclidean'):
40                 distances = [np.linalg.norm(featureset -
41 centroids[centroid]) for centroid in centroids]
42             elif (distance == 'Manhattan'):
43                 distances = [np.linalg.norm(featureset -
44 centroids[centroid], ord=1) for centroid in centroids
45 ]
46             elif (distance == 'Cosine'):
47                 distances = []
48                 for centroid in centroids:
49                     dotP = np.dot(featureset, centroids[
50 centroid].T)
51                     if (isinstance(centroids[centroid],
52 np.float64)):dotP = 0
53                     norm1 = np.linalg.norm(featureset)
54                     if (np.isnan(norm1).any()): norm1 =
55 10**10
56                     norm2 = np.linalg.norm(centroids[
57 centroid])
58                     if (np.isnan(norm2).any()): norm2 =
```

```

10**10
52         distances.append(dotP/(norm2*norm1))
53
54
55
56         minn=np.amin(distances)
57         classification = distances.index(minn)
58         classifications[classification].append(
featureset)
59
60         prev_centroids = dict(centroids)
61
62         for classification in classifications:
63             centroids[classification] = np.average(
classifications[classification],axis=0)
64
65         optimized = True
66
67         for c in centroids:
68             original_centroid = prev_centroids[c]
69             current_centroid = centroids[c]
70             if np.sum((current_centroid-original_centroid
)/original_centroid*100.0) > 0.001:
71                 optimized = False
72
73         if optimized:
74             break
75
76 def predict(data):
77     classification = []
78     for d in data:
79         distances = [np.linalg.norm(d-centroids[centroid
]) for centroid in centroids]
80         classification.append(distances.index(min(
distances)))
81     return classification
82
83 def purity(y_true,y_pred):
84     np.seterr(divide='ignore', invalid='ignore')
85     cm = confusion_matrix(y_true, y_pred, labels=[i
for i in range(10)])
86     return np.sum(np.amax(cm,axis=0)) / np.sum(cm)
87
88 def f_measure(y_true,y_pred):
89     np.seterr(divide='ignore', invalid='ignore')
90     cm = confusion_matrix(y_true, y_pred, labels=[i for i
in range(10)])
91     tp = np.diag(cm)
92     fp=cm.sum(axis=0) - np.diag(cm)
93     fn = cm.sum(axis=1) - np.diag(cm)
94     tn=cm.sum()-(tp+fn+fp)
95     recall = tp/(tp+fn)
96     precision = tp/(tp+fp)
97     f1 = 2*((precision*recall)/(precision+recall))
98     f1[np.isnan(f1)] = 0
99     return f1.sum()
100 print ("-----Kmeans clustering
algorithm-----")
101 print("#####R1-Data Representation
#####")
102 fit(x_train, "Euclidean")

```

```

103 classes = predict(y_test)
104 print("-----Euclidean-----")
105 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(classes))*100))
106 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(classes))*100))
107 fit(x_train, "Manhattan")
108 classes = predict(y_test)
109 print("-----Manhattan-----")
110 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(classes))*100))
111 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(classes))*100))
112 fit(x_train, "Cosine")
113 classes = predict(y_test)
114 print("-----Cosine-----")
115 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(classes))*100))
116 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(classes))*100))
117
118 print("#####R2-Data Representation
    #####")
119
120 histArray=Histogram.histogram(x_train,32)
121 testhistArray=Histogram.histogram(x_test,32)
122
123 fit(histArray, "Euclidean")
124 clusters=predict(testhistArray)
125 print("-----Euclidean-----")
126 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(clusters))*100))
127 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(clusters))*100))
128
129 fit(histArray, "Manhattan")
130 clusters=predict(testhistArray)
131 print("-----Manhattan-----")
132 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(clusters))*100))
133 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(clusters))*100))
134
135 fit(histArray, "Cosine")
136 clusters=predict(testhistArray)
137 print("-----Cosine-----")
138 print("Purity score {:.3f}%".format(purity(y_test,np.
    array(clusters))*100))
139 print("F-measure score {:.3f}%".format(f_measure(y_true=
    y_test, y_pred=np.array(clusters))*100))

```

Αρχικά, αρχικοποιούμε τα κέντρα επιλέγοντας τις πρώτες  $k$  εικόνες από το σύνολο δεδομένων. Στη συνέχεια, για κάθε πρότυπο, βρίσκουμε ποιο κέντρο βρίσκεται πιο κόντα σε αυτό το πρότυπο με βάση μία από τις τρεις συναρτήσεις απόστασης και το τοποθετούμε στη συστάδα του κέντρου αυτού. Έπειτα, ενημερώνουμε την τιμή κάθε κέντρου, θέτοντάς το ίσο με τη μέση τιμή των προτύπων που ανήκουν στη συστάδα του κέντρου αυτού, επαναλαμβάνοντας αυτή τη διαδικασία έως ότου να μην υπάρχει αλλαγή στα κέντρα  $C_k$ . Τα αποτελέσματα της μεθόδου K-means

για τις ζητούμενες αποστάσεις και τις ζητούμενες αναπαραστάσεις των δεδομένων  
είναι τα εξής:

```

-----Kmeans clustering algorithm-----
#####R1-Data Repesantation#####
-----Euclidean-----
Purity score 20.000%
F-measure score 20.000%
-----Manhattan-----
Purity score 20.000%
F-measure score 20.000%
-----Cosine-----
Purity score 10.000%
F-measure score 18.182%
#####R2-Data Repesantation#####
-----Euclidean-----
Purity score 31.020%
F-measure score 75.881%
-----Manhattan-----
Purity score 27.990%
F-measure score 80.428%
-----Cosine-----
Purity score 10.000%
F-measure score 18.182%

```

## 2 Hierarchical clustering

```
1 import math
2 from sklearn.metrics import confusion_matrix
3 import numpy as np
4 import idx2numpy as idx2numpy
5 import Histogram, ReduceDataSet
6
7
8 class Hierarchical:
9     def __init__(self,data):
10         '''
11         constructor of the class, it takes the main data
12         frame as input
13         '''
14         self.data = data
15         self.n_samples, self.n_features = data.shape
16
17     def DistanceMatrix(self,data,distance):
18         '''
19         argument
20         -----
21         data - the dataset whose Similarity matrix we are
22         going to calculate
23         returns
24         -----
25         the distance matrix.
26         '''
27         N = len(data)
28         if (distance == "Cosine"):
29             similarity_mat = np.ones([N, N]) #for cosine np.
30             ones
31         else:
32             similarity_mat = np.zeros([N, N]) #for cosine np.
33             ones
34         for i in range(N):
35             for j in range(N):
36                 similarity_mat[i][j]=self.SimilarityMeasure(data[
37                     i],data[j],distance)
38         return similarity_mat
39
40     def SimilarityMeasure(self,data1, data2, distance):
41         '''
42         arguments
43         -----
44         data1, data2 - vectors, between which we are going to
45         calculate similarity
46         returns
47         -----
48         distance between the two vectors
49         '''
50         N=self.n_features
51         if distance=='Euclidean':
52             # for L2 norm or pythagorean distance
53             dist = 0
54             for i in range(N):
55                 dist += (data1[i] - data2[i])**2
56
57             dist = math.sqrt(dist)
58             return dist
```

```

53
54 if distance=='Cosine':
55     # form cosine similarity = a.b/|a|.|b|
56     dot_prod = 0
57     data1_mod = np.linalg.norm(data1)
58     data2_mod = np.linalg.norm(data2)
59     for x in range(N):
60         dot_prod += data1[x]*data2[x]
61     return (1-dot_prod/(data1_mod*data2_mod))
62
63 if distance=='Manhattan':
64     # for L1 norm or pythagorean distance
65     dist = 0
66     for i in range(N):
67         dist += abs(data1[i] - data2[i] )
68     return dist
69
70 def fit(self,k,distance):
71     '''
72     this method uses the main Divisive Analysis algorithm
73     to do the clustering
74     arguments
75     -----
76     k - integer
77         number of clusters we want
78
79     returns
80     -----
81     cluster_labels - numpy array
82         an array where cluster number of a sample
83         corresponding to
84         the same index is stored
85     '''
86     similarity_matrix = self.DistanceMatrix(self.data,
87     distance) # similarity matrix of the data
88     clusters = [list(range(self.n_samples))] # list
89     of clusters, initially the whole dataset is a single
90     cluster
91     while True:
92         c_diameters = [np.max(similarity_matrix[cluster][:,
93         cluster]) if (len(similarity_matrix[cluster][:,
94         cluster]) != 0) else -1 for cluster in clusters] #
95         cluster diameters
96         max_cluster_dia = np.argmax(c_diameters) #maximum
97         cluster diameter
98         max_difference_index = np.argmax(np.mean(
99         similarity_matrix[clusters[max_cluster_dia]][:,
100         clusters[max_cluster_dia]], axis=1))
101         splinters = [clusters[max_cluster_dia][
102         max_difference_index]] #spinter group
103         last_clusters = clusters[max_cluster_dia]
104         del last_clusters[max_difference_index]
105         while True:
106             split = False
107             for j in range(len(last_clusters))[:-1]:
108                 splinter_distances = similarity_matrix[
109                 last_clusters[j], splinters]
110                 last_distances = similarity_matrix[
111                 last_clusters[j], np.delete(last_clusters, j, axis=0)
112                 ]
113                 if np.mean(splinter_distances) <= np.mean(

```



```

last_distances):
    splinters.append(last_clusters[j])
    del last_clusters[j]
    split = True
    break
    if split == False:
        break
    del clusters[max_cluster_dia]
    clusters.append(splinters)
    clusters.append(last_clusters)
    if len(clusters) == k:
        break

cluster_labels = np.zeros(self.n_samples)
for i in range(len(clusters)):
    cluster_labels[clusters[i]] = i

return (clusters, cluster_labels)

def purity(self, y_true, y_pred):
    np.seterr(divide='ignore', invalid='ignore')
    cm = confusion_matrix(y_true, y_pred, labels=[i for i
        in range(10)])
    return np.sum(np.amax(cm, axis=0)) / np.sum(cm)

def f_measure(self, y_true, y_pred):
    np.seterr(divide='ignore', invalid='ignore')
    cm = confusion_matrix(y_true, y_pred, labels=[i for i
        in range(10)])
    tp = np.diag(cm)
    fp = cm.sum(axis=0) - np.diag(cm)
    fn = cm.sum(axis=1) - np.diag(cm)
    tn = cm.sum() - (tp + fn + fp)
    recall = tp / (tp + fn)
    precision = tp / (tp + fp)
    f1 = 2 * ((precision * recall) / (precision + recall))
    f1[np.isnan(f1)] = 0
    return f1.sum()

trainImagePath = "../dataset/train-images-idx3-ubyte"
trainLabelPath = "../dataset/train-labels-idx1-ubyte"
testImagePath = "../dataset/t10k-images-idx3-ubyte"
testLabelPath = "../dataset/t10k-labels-idx1-ubyte"

x_train = idx2numpy.convert_from_file(trainImagePath)
x_test = idx2numpy.convert_from_file(testImagePath)

y_train = idx2numpy.convert_from_file(trainLabelPath)
y_test = idx2numpy.convert_from_file(testLabelPath)

x_train = x_train.reshape((x_train.shape[0], 28*28))
x_test = x_test.reshape((x_test.shape[0], 28*28))

Hx_train, Hy_train = ReduceDataSet.reduce(x_train, y_train,
    10, 5000)

```

```

156
157 histArray=Histogram.histogram(Hx_train,16)
158
159 print ("-----Hierarchical clustering
        algorithm-----")
160
161 print("#####R2-Data Representation
        #####")
162 Hmodel = Hierarchical(histArray)
163
164 print("-----Euclidean-----")
165 Hclusters, Hlabels = Hmodel.fit(10, "Euclidean")
166 print("Purity score {:.3f}%".format(Hmodel.purity(y_true=
        Hy_train, y_pred=Hlabels)*100))
167 print("F-measure score {:.3f}%".format(Hmodel.f_measure(
        y_true=Hy_train, y_pred=Hlabels)*100))
168
169 print("-----Manhattan-----")
170 Hclusters, Hlabels = Hmodel.fit(10, "Manhattan")
171 print("Purity score {:.3f}%".format(Hmodel.purity(y_true=
        Hy_train, y_pred=Hlabels)*100))
172 print("F-measure score {:.3f}%".format(Hmodel.f_measure(
        y_true=Hy_train, y_pred=Hlabels)*100))
173
174 print("-----Cosine-----")
175 Hclusters, Hlabels = Hmodel.fit(10, "Cosine")
176 print("Purity score {:.3f}%".format(Hmodel.purity(y_true=
        Hy_train, y_pred=Hlabels)*100))
177 print("F-measure score {:.3f}%".format(Hmodel.f_measure(
        y_true=Hy_train, y_pred=Hlabels)*100))
178
179 print("#####R1-Data Representation
        #####")
180
181
182 x_train, y_train = ReduceDataSet.reduce(x_train,y_train
        ,10,2000)
183 x_train=x_train/255.0
184
185 model = Hierarchical(x_train)
186
187 print("-----Euclidean-----")
188 clusters, labels = model.fit(10, "Euclidean")
189 print("Purity score {:.3f}%".format(model.purity(y_true=
        y_train, y_pred=labels)*100))
190 print("F-measure score {:.3f}%".format(model.f_measure(
        y_true=y_train, y_pred=labels)*100))
191
192 print("-----Manhattan-----")
193 clusters, labels = model.fit(10, "Manhattan")
194 print("Purity score {:.3f}%".format(model.purity(y_true=
        y_train, y_pred=labels)*100))
195 print("F-measure score {:.3f}%".format(model.f_measure(
        y_true=y_train, y_pred=labels)*100))
196
197 print("-----Cosine-----")
198 clusters, labels = model.fit(10, "Cosine")
199 print("Purity score {:.3f}%".format(model.purity(y_true=
        y_train, y_pred=labels)*100))
200 print("F-measure score {:.3f}%".format(model.f_measure(
        y_true=y_train, y_pred=labels)*100))

```

Αρχικά, ξεκινάμε με όλα τα σημεία σε μία συστάδα, η οποία στη συνέχεια, χωρίζεται σε δύο λιγότερο όμοιες συστάδες. Υπολογίζουμε τον distance matrix  $C(i, j)$ , ο οποίος αναπαριστά τις αποστάσεις των σημείων  $i, j$  (σύμφωνα με κάποια από τις τρεις συναρτήσεις απόστασης.) Έπειτα, η συστάδα διασπάται με βάση τη μεγαλύτερη μέση απόσταση σε δύο νέες συστάδες και τα εναπομείναντα σημεία μετακινούνται, εάν η διαφορά μεταξύ της μέσης απόστασης και της απόστασης με την άλλη συστάδα είναι θετική. Στη συνέχεια, υπολογίζουμε τη διάμετρο της κάθε συστάδας, η οποία ορίζεται ως η απόσταση μεταξύ δύο οποιονδήποτε σημείων της συστάδας και επιλέγεται η συστάδα με τη μεγαλύτερη διάμετρο για περαιτέρω διάσπαση. Η παραπάνω διαδικασία, επαναλαμβάνεται έως ότου πάρουμε τον επιθυμητό αριθμό συστάδων. Τα αποτελέσματα του Hierarchical clustering για τις ζητούμενες αποστάσεις και τις ζητούμενες αναπαραστάσεις των δεδομένων είναι τα εξής:

```

-----Hierarchical clustering algorithm-----
#####R2-Data Representantation#####
-----Euclidean-----
Purity score 28.240%
F-measure score 12.355%
-----Manhattan-----
Purity score 26.260%
F-measure score 42.998%
-----Cosine-----
Purity score 29.480%
F-measure score 74.590%
#####R1-Data Representantation#####
-----Euclidean-----
Purity score 47.400%
F-measure score 79.292%
-----Manhattan-----
Purity score 43.300%
F-measure score 98.818%
-----Cosine-----
Purity score 40.150%
F-measure score 98.405%

```

### 3 K-medoids algorithm

```
1 import numpy as np
2 from math import log2
3 from sklearn.metrics import confusion_matrix
4 import idx2numpy as idx2numpy
5 import Histogram, ReduceDataSet
6
7
8 class k_medoids:
9     def __init__(self, k=10): # total number of clusters
10         =10
11         self.k = k
12
13     def fit(self, data):
14
15         # For 1st clustering
16         self.medoids = []
17         self.cost = 0
18         self.clusters = []
19         for i in range(self.k):
20             self.medoids.append(data[i]) # selecting k
21             random points out of the dataset as the medoids (1st
22             step)
23         for j in range(self.k):
24             self.clusters.append([])
25
26         for point in data:
27             # list containing distances of each point in
28             the dataset from the medoids and finding minimum
29             distance
30
31             distance = [np.sum(np.where(point != 0, point
32             * np.log(point / m), 0)) for m in self.medoids]
33             min_distance = min(distance)
34             # calculating the cost (cost is the distance
35             of each point from its medoid i.e., total of min
36             distances)
37             self.cost += min_distance
38             o = distance.index(min_distance)
39             # clustering (i.e., associating each point to
40             the closest medoid)
41             self.clusters[o].append(point)
42
43         # now, at each iteration we are finding the
44         replacement of our medoids
45         # we will stop if there is no changes in the cost
46         while (True):
47
48             # lists that will store new clusters, new
49             cost and the new medoids
50             new_clusters = []
51             self.new_cost = 0
52             self.new_medoids = []
53             for j in range(self.k):
54                 new_clusters.append([])
55
56             for j in range(self.k):
57                 # list for storing the total distance of
58                 a particular point from all other points in the same
```

```

47         cluster
48         dist_with_each_point_in_same_cluster = []
49         l = [] # list for storing the distance
50         of each and every point from all other points of the
51         same cluster
52         for point in self.clusters[j]:
53             dist_with_each_point_in_same_cluster
54             = sum([np.sum(np.where(point != 0, point * np.log(
55                 point / d), 0)) for d in self.clusters[j]])
56             l.append(
57                 dist_with_each_point_in_same_cluster)
58             minima = min(l) # finding the minimum
59             distance
60             q = l.index(minima)
61             # point with the min distance from all
62             the points in the same cluster is taken as the new
63             medoid
64             self.new_medoids.append(self.clusters[j][
65                 q])
66
67         # now, finding the new clustering and the
68         new cost
69         for point in data:
70             # list containing distances of each point
71             in the dataset from the new medoids and finding
72             minimum distance
73             distance = [np.sum(np.where(point != 0,
74                 point * np.log(point / m), 0)) for m in self.
75                 new_medoids]
76             min_distance = min(distance)
77             # calculating the new cost
78             self.new_cost += min_distance
79             o = distance.index(min_distance)
80             # new clustering
81             new_clusters[o].append(point)
82
83         # if the cost decreases with the new medoids
84         we will replace the old medoids, old clusters and the
85         old cost with the new one
86         # if the cost increases with the new medoids
87         we will not replace the old medoids, old clusters and
88         the oldcost with the new one
89
90         if self.new_cost < self.cost:
91             self.medoids = self.new_medoids
92             self.clusters = new_clusters
93             self.cost = self.new_cost
94
95         # if the cost remains same we will stop and
96         come out of the iterations
97
98         elif abs(self.new_cost - self.cost) <= 0.01:
99             break
100
101     # final clustering according to the new medoids
102     def predict(self, test_data):
103         pred = []
104         for point in test_data:
105             # finding distance of each point from the

```

```

final medoids and get the min distance
88     distance = [np.sum(np.where(point != 0, point
    * np.log(point / m), 0)) for m in self.medoids]
89     minDistance = min(distance)
90     l = distance.index(minDistance)
91     pred.append(l) # associating each point to
the closest medoid
92     return pred
93
94 def purity(self, y_true, y_pred):
95     np.seterr(divide='ignore', invalid='ignore')
96     cm = confusion_matrix(y_true, y_pred, labels=[i
for i in range(10)])
97     return np.sum(np.amax(cm, axis=0)) / np.sum(cm)
98
99 def f_measure(self, y_true, y_pred):
100     np.seterr(divide='ignore', invalid='ignore')
101     cm = confusion_matrix(y_true, y_pred, labels=[i
for i in range(10)])
102     tp = np.diag(cm)
103     fp = cm.sum(axis=0) - np.diag(cm)
104     fn = cm.sum(axis=1) - np.diag(cm)
105     tn = cm.sum() - (tp + fn + fp)
106     recall = tp / (tp + fn)
107     precision = tp / (tp + fp)
108     f1 = 2 * ((precision * recall) / (precision +
recall))
109     f1[np.isnan(f1)] = 0
110     return f1.sum()
111
112
113 trainImagePath = "../dataset/train-images-idx3-ubyte"
114 trainLabelPath = "../dataset/train-labels-idx1-ubyte"
115 testImagePath = "../dataset/t10k-images-idx3-ubyte"
116 testLabelPath = "../dataset/t10k-labels-idx1-ubyte"
117
118 x_train = idx2numpy.convert_from_file(trainImagePath)
119 x_test = idx2numpy.convert_from_file(testImagePath)
120
121 y_train = idx2numpy.convert_from_file(trainLabelPath)
122 y_test = idx2numpy.convert_from_file(testLabelPath)
123
124 x_train = x_train.reshape((x_train.shape[0], 28 * 28))
125 x_test = x_test.reshape((x_test.shape[0], 28 * 28))
126
127 x_train,y_train = ReduceDataSet.reduce(x_train,y_train
,10,6000)
128 print(x_train.shape)
129 bins = 16
130 histArray = Histogram.histogram(x_train,bins)
131 testhistArray=Histogram.histogram(x_test,bins)
132
133 model = k_medoids()
134 model.fit(histArray)
135 clusters = model.predict(testhistArray)
136
137 print ("-----Kmedoid clustering
algorithm-----")
138 print("#####R2-Data Representation
#####")
139 print("Purity score {:.3f}%".format(model.purity(y_true=

```

```

        y_test, y_pred=clusters)*100))
140 print("F-measure score {:.3f}%".format(model.f_measure(
        y_true=y_test, y_pred=clusters)*100))

```

Αρχικά, αρχικοποιούμε τις διαμέσους (medoids) επιλέγοντας τις πρώτες  $k$  εικόνες από το σύνολο δεδομένων. Στη συνέχεια, για κάθε πρότυπο, βρίσκουμε ποιο medoid είναι περισσότερο όμοιο σε αυτό το πρότυπο με βάση την symmetric Kullback-Libler distance και το τοποθετούμε στη συστάδα του medoid αυτού. Έπειτα, ενημερώνουμε την τιμή κάθε medoid, θέτοντάς το ίσο με το σημείο του cluster που έχει τη μεγαλύτερη ομοιότητα με όλα τα υπόλοιπα σημεία, επαναλαμβάνοντας αυτή τη διαδικασία έως ότου να μην υπάρχει αλλαγή στα κέντρα  $C_k$ . Τα αποτελέσματα της μεθόδου K-medoids είναι τα εξής:

---

```

-----Kmedoid clustering algorithm-----
#####R2-Data Represantation#####
Purity score 17.880%
F-measure score 89.048%

```

## 4 Βοηθητική συνάρτηση histogram

```
1 import numpy as np
2 def histogram(X,bins):
3     histArray=[]
4     for i in X:
5         histArray.append(np.array(np.histogram(i,bins)[0]))
6     return np.array(histArray, dtype='float64')
```

Η συνάρτηση histogram που βρίσκεται στο αρχείο Histogram.py παίρνει ως όρισμα ένα dataset και τον αριθμό των bins που επιθυμούμε και δημιουργεί το τελικό dataset που περιέχει το ιστόγραμμα της κάθε εικόνας.

## 5 Βοηθητική συνάρτηση reduce

```
1 import numpy as np
2 def reduce(X,Y,k, n):
3     dataset = []
4     labels = []
5     for i in range(k):
6         c=n//k
7         j=0
8         while c>0:
9             if Y[j] == i:
10                 dataset.append(X[j])
11                 labels.append(Y[j])
12                 c-=1
13             j+=1
14     return np.array(dataset), np.array(labels)
```

Η συνάρτηση reduce που βρίσκεται στο αρχείο ReduceDataSet.py παίρνει ως όρισμα ένα dataset μαζί με τα labels, το πλήθος των κλάσεων και τέλος τον αριθμό των στοιχείων που θα έχει το μειωμένο (τελικό) dataset. Για κάθε κατηγορία, υπολογίζουμε το πλήθος των στοιχείων που θα έχει αυτή η κατηγορία και έπειτα τα τοποθετούμε στο τελικό dataset.

## 6 Η βέλτιστη μέθοδος

Έπειτα από σύγκριση των μέτρων αξιολόγησης Purity, F-measure καταλήγουμε πως η βέλτιστη μέθοδος για την  $R_1$  αναπαράσταση δεδομένων είναι η μέθοδος της ιεραρχικής συσταδοποίησης με χρήση της ευκλείδιας απόστασης καθώς έχει το μεγαλύτερο Purity (47.4%) και για την  $R_2$  αναπαράσταση δεδομένων είναι η μέθοδος της ιεραρχικής συσταδοποίησης με χρήση της συνημιτονοειδούς απόστασης καθώς έχει το μεγαλύτερο Purity (29.48%).