

## Program -3 DBSCAN

**Dataset:** – This dataset contains 8580 text records in sparse format. In this dataset labels are not provided.

**Objective :** The objective of this assignment is to use Kmeans and variation of DBSCAN to assign each text record to a cluster ranging from 1 to  $k > 100$

### Approach –

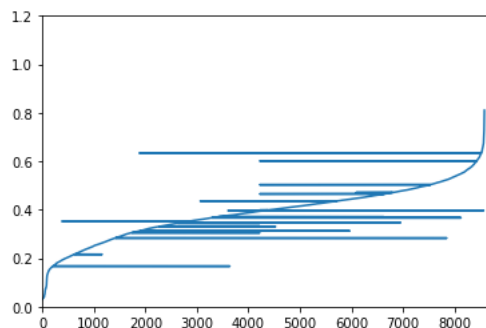
1) All the records are read using **csr\_read** function and was saved into mat1 variable

2) mat1 is transformed in to a count matrix to a normalized tf or tf-idf representation using **TfidfTransformer().fit\_transform** module. Library to **import from sklearn.feature\_extraction.text import TfidfTransformer**

3) Performed Dimensionality Reduction using SVD by keeping the following parameters **n\_components=200, n\_iter=7, random\_state=42**. Resultant matrix is saved into **mat\_dr**.

4) After that each value is normalized to unit scale using **normalizer.fit\_transform** and saved into **points**. Also, a copy of points is made as **points\_or** for future reference. Variable **explained\_variance** gives the Explained Variance of SVD and I got it as **61 Percent**.

5) Epsilon value is computed using K-nearest neighbor and graph is plotted for the same. Parameter used **n\_neighbors = 30, metric='cosine'**



6) This matrix saved is then given to KMeans to find clusters . Number of clusters used is 500. After that centroids obtained for all the cluster is saved in **Cluster\_center** and labels for the respective clusters are saved in **cluster\_label**.

7) In order to make cluster of clusters I decided to use center points obtained through Kmeans and find core points , border points and noise points for it. I tried with different values of eps and min points but combination of **eps = .60** and **minpts = 3** worked best for me. I got **282 core points** with this , **90 border points** , and **128 Noise points**

8) `G = nx.Graph()` was used to Invoke graph instance to visualize the cluster

9) Number of clusters obtained were 7

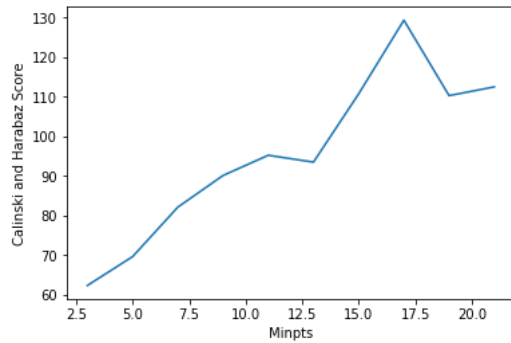
10) In order to make the prediction cluster\_label is again copied into variable named KmeansLabels. A dictionary names cluster\_dict was made and all the labels obtained through kmeans was converted into the corresponding DBSCAN labels . Also an additional cluster was made to adjust all the noise points. Final prediction was saved into clustpred file.

11) For the purpose of inter cluster evaluation I have used calinski harabaz metrics. I ran the code for different epsilon and k values and got the following table. I plotted the graph for epsilon value 0.7.

Rank. -15

NMI – 0.4781

<b>k\eps</b>	<b>0.7</b>
<b>3</b>	<b>62.346865</b>
<b>5</b>	<b>69.6394353</b>
<b>7</b>	<b>82.1397973</b>
<b>9</b>	<b>90.1328761</b>
<b>11</b>	<b>95.2451444</b>
<b>13</b>	<b>93.5108298</b>
<b>15</b>	<b>110.755648</b>
<b>17</b>	<b>129.380118</b>
<b>19</b>	<b>110.314061</b>
<b>21</b>	<b>112.519327</b>



```
import numpy as np
import pandas as pd
import random
from collections import defaultdict
from scipy.sparse import csr_matrix
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import Normalizer
from sklearn.neighbors import NearestNeighbors
import matplotlib.pyplot as plt
import numpy
```

```
# In[849]:
```

```
#reading csr file and returning csr matrix by calculating ind,ptr and values in the file.
```

```
def csr_read(fname, ftype="csr", nidx=1):

    with open(fname) as f:
        lines = f.readlines()
        nrows = len(lines)
        ncols = 0
        nnz = 0
        for i in range(nrows):
            p = lines[i].split()
            if len(p) % 2 != 0:
                raise ValueError("Invalid CSR matrix")
            nnz += len(p)/2
            nnz=int(nnz)
            for j in range(0, len(p), 2):
                cid = int(p[j]) - nidx
                if cid+1 > ncols:
                    ncols = cid+1

        val = np.zeros(nnz, dtype=np.float)
        ind = np.zeros(nnz, dtype=np.int)
        ptr = np.zeros(nrows+1, dtype=np.long)
        n = 0
        for i in range(nrows):
            p = lines[i].split()
            for j in range(0, len(p), 2):
                ind[n] = int(p[j]) - nidx
                val[n] = float(p[j+1])
                n += 1
            ptr[i+1] = n

        assert(n == nnz)

    return csr_matrix((val, ind, ptr), shape=(nrows, ncols), dtype=np.float)
```

```

5
6 #performing dimensionality reduction using SVD to 150 components and normalising the output to
  inhibit impact of document length on distance computation
7 print("Performing dimensionality reduction using LSA")
8 svd = TruncatedSVD(n_components=200,n_iter=7, random_state=42)
9 normalizer = Normalizer(copy=False)
10 mat_dr=svd.fit_transform(mat1)
11 points=normalizer.fit_transform(mat_dr)
12 points_or=points
13 explained_variance = svd.explained_variance_ratio_.sum()
14 print(explained_variance)
15
16
17 # In[854]:
18
19
20 #Computing eps for minpts using K-distance graph and elbow point in it.
21 nbrs = NearestNeighbors(n_neighbors=30,metric='cosine').fit(points)
22 distances, indices = nbrs.kneighbors(points)
23 t=distances[:,~1]
24 i=indices[:,0]
25 a=np.sort(t)
26 plt.axis([0, 8680, 0, 1.2])
27 plt.plot(i,a)
28
29
30 # In[855]:
31
32
33 import sklearn
34 from sklearn.cluster import KMeans
35 from mpl_toolkits.mplot3d import Axes3D
36 from sklearn.preprocessing import scale
37
38 import sklearn.metrics as sm
39 from sklearn.metrics import confusion_matrix,classification_report
40
41
42 # In[856]:
43
44
45 clustering = KMeans(n_clusters = 500, random_state = 42)
46 clustering.fit(points)
47
48
49 # In[857]:
50
51
52 print(clustering.labels_)
53 cluster_label= clustering.labels_
54 # cluster_len
55 #cluster_len
56 X1=pd.Series(clustering.labels_)
57 X1.value_counts()
58
59
60 # In[858]:
61
62
63 #Plotting the data points again on the graph and visualize how the data has been clustered
64 plt.scatter(mat_dr[:,0],mat_dr[:,1], c=clustering.labels_, cmap='rainbow')
65
66
67 # In[859]:
68
69
70 #Plotting the data points again on the graph and visualize the centroids in black
71 plt.scatter(mat_dr[:,0], mat_dr[:,1], c=clustering.labels_, cmap='rainbow')
72 plt.scatter(clustering.cluster_centers_[0],clustering.cluster_centers_[1], color='black')
73
74
75 # In[860]:
76
77
78 Cluster_center = np.array(clustering.cluster_centers_)
79
80
81 # In[861]:
82
83
84 Cluster_center.shape
85

```

```

1 #DBSCAN tryout using networkx
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import networkx as nx
5
6
7 # In[880]:
8
9
0 points = Cluster_center
1 eps = .60
2 minPts = 3
3
4
5 # In[881]:
6
7
8 # Find core points with minPts as 5
9 neighborhoods = []
0 core = []
1 border = []
2 noise = []
3
4 for i in range(len(points)):
5     neighbors = []
6     for p in range(0, len(points)):
7         # If the distance is below eps, p is a neighbor
8         if np.linalg.norm(points[i] - points[p]) < eps:
9             neighbors.append(p)
0     neighborhoods.append(neighbors)
1     # If neighborhood has at least minPts, i is a core point
2     if len(neighbors) > minPts :
3         core.append(i)
4
5 print("core: ", core)
6
7
8 # In[882]:
9
0
1 print(len(core))
2
3
4 # In[883]:
5
6
7 # Find border points
8 for i in range(len(points)):
9     neighbors = neighborhoods[i]
0     # Look at points that are not core points
1     if len(neighbors) <= minPts:
2         for j in range(len(neighbors)):
3             # If one of its neighbors is a core, it is also in the core point's neighborhood,
4             # thus it is a border point rather than a noise point
5             if neighbors[j] in core:
6                 border.append(i)
7                 # Need at least one core point...
8                 break
9
0 print("border: ", border)
1 #print(len(border))
2
3
4 # In[889]:
5
6
7 print(len(border))
8
9
0 # In[884]:
1
2
3 #Find noise points
4 for i in range(len(points)):
5     if i not in core and i not in border:
6         noise.append(i)
7
8 print("noise", noise)
9 #print("Length of noise points :", len(noise))
0
1
2 # In[890]:
3
4
5 print("Length of noise points :", len(noise))

```

---

```

279
280
281 # Invoke graph instance to visualize the cluster
282 G = nx.Graph()
283
284
285 # In[892]:
286
287
288 # Add nodes -- core points + border points
289 nodes = core+border
290 G.add_nodes_from(nodes)
291
292
293 # In[893]:
294
295
296 # Create neighborhood
297 for i in range(len(nodes)):
298     for p in range(len(nodes)):
299         # If the distance is below the threshold, add a link in the graph.
300         if p != i and np.linalg.norm(points[nodes[i]] - points[nodes[p]]) <= eps:
301             G.add_edges_from([(nodes[i], nodes[p])])
302
303
304 # In[894]:
305
306
307 # List the connected components / clusters
308 clusters = list(nx.connected_components(G))
309 print("# clusters:", len(clusters))
310 print("clusters: ", clusters)
311
312
313 # In[878]:
314
315
316 # Visualise the graph
317 plt.subplot(111)
318 nx.draw_circular(G, with_labels=True, font_weight='bold')
319 plt.show()
320
321
322 # In[879]:
323
324
325
326 kmeansLabels = cluster_label
327
328 cluster_dict = {}
329 for i, cluster in enumerate(clusters):
330     for cluster_center in cluster:
331         cluster_dict[cluster_center] = i
332 #np.array(clusters)
333
334 for cluster_center in noise:
335     cluster_dict[cluster_center] = -1
336
337 # To convert kmeans labels to dbscan labels
338 labels_final = [cluster_dict[x] for x in kmeansLabels]
339
340 labels_final = np.array(labels_final)
341 labels_final[labels_final==-1] = max(labels_final)+1
342 labels_final[labels_final==0] = max(labels_final)+1
343
344 print(pd.Series(labels_final).value_counts())
345
346
347 # In[895]:
348
349
350 fh = open("clustpred.dat", "w")
351 for x in labels_final:
352     fh.write("{}\n".format(x))
353 fh.close()
354

```

```
from sklearn import metrics
```

```
# In[897]:
```

```
s=metrics.calinski_harabaz_score(points_or, labels_final)
```

```
print(s)
```

```
# In[898]:
```

```
#Precomputed scores
```

```
score=
```

```
[62.34686498,69.6394,82.13979732,90.13287605,95.24514444,93.51082975,110.755648,129.3801181,110.3140608,112.5193268]
```

```
k=[3,5,7,9,11,13,15,17,19,21]
```

```
plt.xlabel('Minpts')
```

```
plt.ylabel('Calinski and Harabaz Score')
```

```
plt.plot(k,score)
```

---