# ▾ EFFECTIVE DENSITY BASED CLUSTERING

This paper provides a solution to perform density-based clustering on datasets that have incomplete data without loss of quality of clusters.

CI clustering with intermediate clustering based on the KD tree. LI clustering: Create a kD tree with the entire dataset and use the information of neighboring points to predict the missing values of a cluster. DBSCAN to achieve a result(cluster).

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.impute import KNNImputer
4 from sklearn import datasets
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.cluster import DBSCAN
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9 from sklearn.model_selection import train_test_split
```

+ Code    + Text

```
1 data=pd.read_csv("synb.csv")
2 print(data.head())
3 print(data.columns)
4 Before_imputation = pd.DataFrame(data)
5 print("Data Before performing imputation\n",Before_imputation.head())
```

```
        x      y  class
0  0.228  0.559      1
1  0.216  0.528      1
2  0.221  0.552      1
3  0.215  0.538      1
4  0.224  0.548      1
Index(['x', 'y', 'class'], dtype='object')
Data Before performing imputation
        x      y  class
0  0.228  0.559      1
1  0.216  0.528      1
2  0.221  0.552      1
3  0.215  0.538      1
4  0.224  0.548      1
```

```
1 from sklearn.impute import SimpleImputer
2 df=pd.DataFrame(data)
3 mean_imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
4
5 # Fit the imputer on to the dataset
6 mean_imputer = mean_imputer.fit(df)
7
8
9 # Apply the imputation
10 results = mean_imputer.transform(df.values)
```

```
11 results
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid featu
  warnings.warn(
array([[0.228, 0.559, 1.   ],
       [0.216, 0.528, 1.   ],
       [0.221, 0.552, 1.   ],
       ...,
       [0.513, 0.233, 2.   ],
       [0.506, 0.221, 2.   ],
       [0.515, 0.26 , 2.   ]])
```

```
 1
 2 imputer = KNNImputer(n_neighbors=2)
 3 After_imputation = imputer.fit_transform(Before_imputation)
 4 ai=pd.DataFrame(After_imputation)
 5 cdata = ai.to_csv("TrainingDataset.csv")
 6 data1=pd.read_csv("TrainingDataset.csv",usecols=(1,2))
 7 data1["class"]=Before_imputation["class"]
 8 data1.head()
 9 print(data1.shape)
10 data1.isnull().any().any()
11 x=data1.loc[:,['0','1']].values
12 print(x.shape)
13 data1
```

```
(4811, 3)
(4811, 2)
```

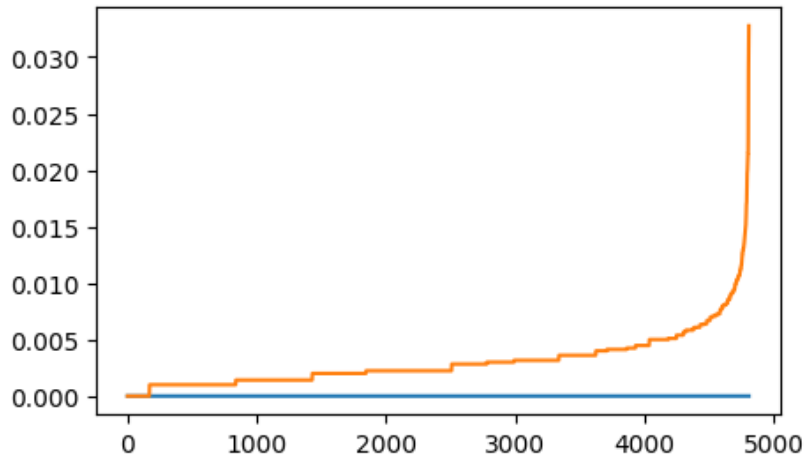|      | 0     | 1     | class |
|------|-------|-------|-------|
| 0    | 0.228 | 0.559 | 1     |
| 1    | 0.216 | 0.528 | 1     |
| 2    | 0.221 | 0.552 | 1     |
| 3    | 0.215 | 0.538 | 1     |
| 4    | 0.224 | 0.548 | 1     |
| ...  | ...   | ...   | ...   |
| 4806 | 0.507 | 0.269 | 2     |
| 4807 | 0.526 | 0.237 | 2     |
| 4808 | 0.513 | 0.233 | 2     |
| 4809 | 0.506 | 0.221 | 2     |
| 4810 | 0.515 | 0.260 | 2     |

4811 rows × 3 columns

```
1 from sklearn.neighbors import NearestNeighbors
2 neighb = NearestNeighbors(n_neighbors=2)
3 nbrs=neighb.fit(x)
```

```
4 distances,indices=nbrs.kneighbors(x)
5 distances = np.sort(distances, axis = 0)
6 distances = distances[:,]
7 plt.rcParams['figure.figsize'] = (5,3)
8 plt.plot(distances)
9 plt.show()
```
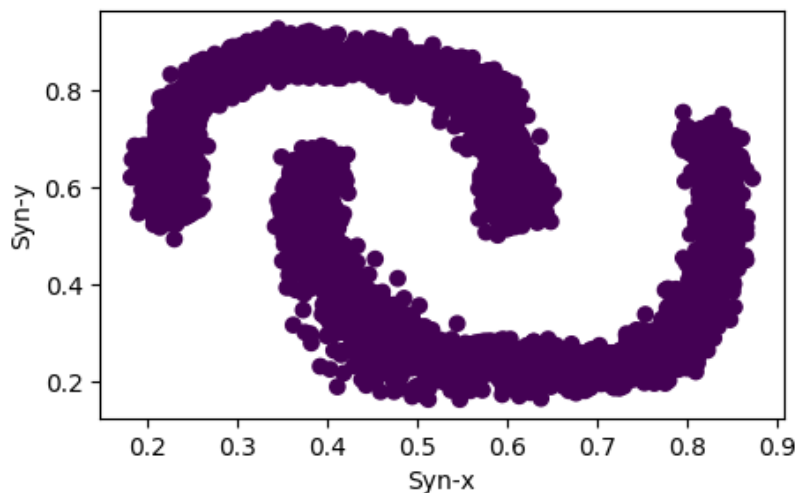


```
1 from sklearn.cluster import DBSCAN
2 dbscan = DBSCAN(eps = 8, min_samples = 4).fit(x)
3 labels = dbscan.labels_
4 plt.scatter(x[:, 0], x[:,1], c = labels, cmap= "viridis")
5 plt.xlabel("Syn-x")
6 plt.ylabel("Syn-y")
7 plt.show()
```



```
1 Before_imputation = pd.DataFrame(data)
2 imputer = KNNImputer(n_neighbors=2)
3 After_imputation = imputer.fit_transform(Before_imputation)
4 ai=pd.DataFrame(After_imputation)
5 cdata = ai.to_csv("new.csv")
6 data1=pd.read_csv("new.csv",usecols=(1,2))
```

```
7 data1["class"]=Before_imputation["class"]
8 data1
```

|      | 0     | 1     | class |
|------|-------|-------|-------|
| 0    | 0.228 | 0.559 | 1     |
| 1    | 0.216 | 0.528 | 1     |
| 2    | 0.221 | 0.552 | 1     |
| 3    | 0.215 | 0.538 | 1     |
| 4    | 0.224 | 0.548 | 1     |
| ...  | ...   | ...   | ...   |
| 4806 | 0.507 | 0.269 | 2     |
| 4807 | 0.526 | 0.237 | 2     |
| 4808 | 0.513 | 0.233 | 2     |
| 4809 | 0.506 | 0.221 | 2     |
| 4810 | 0.515 | 0.260 | 2     |

4811 rows × 3 columns

```
 1 from sklearn.tree import DecisionTreeClassifier
 2 data1.head()
 3 X=data1.drop("class",axis=1)
 4 Y=data1["class"]
 5 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.30, random
 6 dt=DecisionTreeClassifier()
 7 dt.fit(X_train,y_train)
 8 res=dt.score(X_test,y_test)
 9 Xa=data1.drop("class",axis=1)
10 Ya=data1["class"]
11 X_1train, X_1test, y_1train, y_1test = train_test_split(Xa, Ya, test_size=0.30,
12 res
```

    1.0

```
 1 ## KD TREE
 2
 3 class TreeNode(object):
 4     def __init__(self, v1,v2):
 5         self.x = v1
 6         self.y=v2
 7         self.left = None
 8         self.right = None
 9
10 def create_tree(arr):
11     if not arr:
```

```
12          return None
13      mid_num = len(arr)//2
14      node = TreeNode(arr[mid_num][0],arr[mid_num][1])
15      node.left = create_tree(arr[:mid_num])
16      node.right = create_tree(arr[mid_num+1:])
17      return node
18 def RangeSearch(elem,tree,temp,arr,num=0):
19   if not tree or num>=len(elem):
20      return None
21   if not elem:
22      return None
23   # print(elem[num][0],tree.x,tree.y)
24   if elem[num][0]==tree.x:
25     # print("ahsh")
26     if tree.left is None or tree.right is None:
27        arr.append(tree)
28     else:
29        arr.append(tree.left)
30        arr.append(tree.right)
31     num+=1
32     RangeSearch(elem,temp,temp,arr,num)
33   else:
34        RangeSearch(elem,tree.right,temp,arr,num)
35        RangeSearch(elem,tree.left,temp,arr,num)
36
37 def Predict(elem,tree,arr):
38   if not tree:
39      return None
40   if not elem:
41      return None
42   # print(elem[num][0],tree.x,tree.y)
43   if elem[0]==tree.x:
44     if tree.left is None or tree.right is None:
45        arr.append(tree)
46     else:
47        arr.append(tree.left)
48        arr.append(tree.right)
49   else:
50        Predict(elem,tree.right,arr)
51        Predict(elem,tree.left,arr)
52
53 def insert(tree, x,y):
54   if tree.x:
55     if x < tree.x:
56        if tree.left is None:
57            tree.left = TreeNode(x,y)
58        else:
59            insert(tree.left,x,y)
```

```
60    elif x > tree.x:
61      if tree.right is None:
62        tree.right = TreeNode(x,y)
63      else:
64        insert(tree.right,x,y)
65
66 def getlist(tree,f_list):
67   if tree.left:
68       getlist(tree.left,f_list)
69   if tree.right:
70       getlist(tree.right,f_list)
71   f_list.append([tree.x,tree.y])
```

```
 1 ####################
 2 """
 3  CI-Clustering
 4 Input: Xcomplete and Xincomplete
 5 Output: Vector y denoting the clustering result
 6 1: Normalize.Xcomplete/
 7 2: Normalize.Xincomplete/
 8 3: tree KD-TREE.Xcomplete/
 9 4: neighborhoods RangeSearch.Xcomplete; tree/
10 5: clusters getClusters.Xcomplete; neighborhoods/
11 6: for x in Xincomplete do
12 7: x predict.x; tree; clusters/
13 8: insert.tree; x/
14 9: update neighborhoods.neighborhoods/
15 10: update clusters.clusters/
16 11: en
17 """
18 from sklearn.metrics import f1_score
19 # from anytree import Node, RenderTree
20 from sklearn.cluster import DBSCAN
21 from sklearn import preprocessing
22 from sklearn.neighbors import KDTree
23 from sklearn.cluster import DBSCAN
24 complete=X_1train.values.tolist()
25 uncomplete=X_train.values.tolist()
26 complete_norm = preprocessing.normalize(X_1train)
27 uncomplete_norm=preprocessing.normalize(X_train)
28 # print(complete_norm.tolist())
29 # print(X_1train)
30 tree=create_tree(complete_norm.tolist())
31 # print(tree)
32 arr=[]
33 num=0
34 RangeSearch(complete_norm.tolist()[:20],tree,tree,arr,num)
```

```
35 weights=[]
36 for i in arr:
37   weights.append(i.x)
38   weights.append(i.y)
39 # print(arr)
40 dbscan = DBSCAN(eps = 8, min_samples = 4)
41 train=dbscan.fit(complete_norm)
42 clusters=train.labels_
43 # print(0 in clusters.tolist())
44 for x in uncomplete_norm:
45   neigh=[]
46   x=x.reshape(1,-1)
47   # print(x)
48   Predict(x.tolist()[0],tree,neigh)
49   # print(neigh)
50   miss_y=0
51   miss_x=0
52   for i in neigh:
53     # print(i)
54     miss_y+=i.y
55     miss_x+=i.x
56   insert(tree,miss_x,miss_y)
57 f_list=[]
58 getlist(tree,f_list)
59   # print(f_list)
60 res=dt.predict(f_list)
61 print(res)
62 # check=X_1test.values.tolist().extend(X_test.values.tolist())
63 frames = [y_train, y_1train]
64 result = pd.concat(frames)
65 print(result.shape)
66 print(res.shape)
67 si=result.shape[0]-res.shape[0]
68 output=f1_score(res,result[si:],average='weighted')
69 output
70
71
72
73
74 # labels = dbscan.labels_
75 # n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
76 # ans=tree.query_radius(X_train[:1], r=0.3)
77 # appender=np.array(normalized)
78 # dbscan = DBSCAN(eps = 8, min_samples = 4)
79 # # ans
80 # for x in nomr2:
81 #   x=x.reshape(1,-1)
82 #   val=dbscan.fit_predict(x)
```

```
83 #    np.insert(appender,obj=len(appender),values=val)
84 #    # tree=array_to_bst(appender)
85 #    labels = dbscan.labels_
86 #    # n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
87 # appender
88
```

```
[1 1 1 ... 1 2 2]
(6734,)
(6669,)
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid featu
  warnings.warn(
0.5032554188832492
```

```
 1 #LI Clustering
 2 from sklearn.metrics import f1_score
 3 # from anytree import Node, RenderTree
 4 from sklearn.cluster import DBSCAN
 5 from sklearn import preprocessing
 6 from sklearn.neighbors import KDTree
 7 complete=X_1train.values.tolist()
 8 uncomplete=X_train.values.tolist()
 9 complete_norm = preprocessing.normalize(X_1train)
10 uncomplete_norm=preprocessing.normalize(X_train)
11 # print(complete_norm.tolist())
12 # print(X_1train)
13 tree=create_tree(complete_norm.tolist())
14 dbscan = DBSCAN(eps = 8, min_samples = 4)
15 train=dbscan.fit(complete_norm)
16 clusters=train.labels_
17 # print(0 in clusters.tolist())
18 for x in uncomplete_norm:
19   neigh=[]
20   x=x.reshape(1,-1)
21   # print(x)
22   Predict(x.tolist()[0],tree,neigh)
23   # print(neigh)
24   miss_y=0
25   miss_x=0
26   for i in neigh:
27     # print(i)
28     miss_y+=i.y
29     miss_x+=i.x
30   insert(tree,miss_x,miss_y)
31 f_list=[]
32 getlist(tree,f_list)
33 train=dbscan.fit(f_list)
34 clusters=train.labels_
```

```
35 res=dt.predict(f_list)
36 print(res)
37
```

```
[1 1 1 ... 1 2 2]
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid featu
  warnings.warn(
```
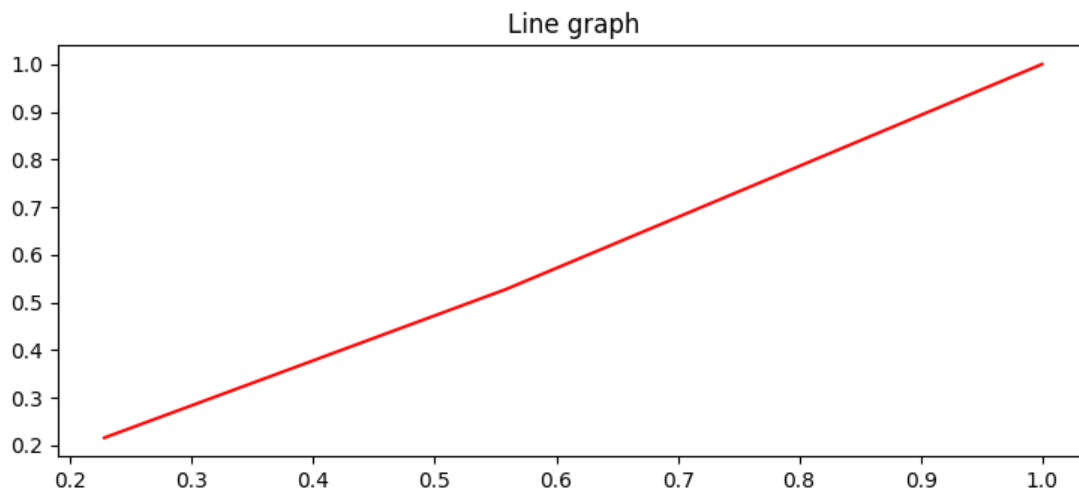
```
1
2 f_list
```

```
[[0.27078309672430856, 0.962640386919432],
 [0.2683244075628877, 0.9633286107585642],
 [0.27335843910850066, 0.9619122432780259],
 [0.2714601650218062, 0.9624496759864037],
 [0.2755596101787932, 0.9612839857389238],
 [0.27506063191739133, 0.961426881655181],
 [0.27999999999999997, 0.96],
 [0.2791280434473638, 0.960253891094041],
 [0.2809958715857758, 0.9597089768006498],
 [0.2778914668170219, 0.9606124778860017],
 [0.2823263743670628, 0.9593184134252554],
 [0.28229584288523907, 0.959327398279499],
 [0.28202020961790236, 0.9594084642982228],
 [0.28129180444881063, 0.9596222802488134],
 [0.28309107310433124, 0.959093032155191],
 [0.2844470071718319, 0.9586917649124705],
 [0.2843842936854435, 0.9587103699788749],
 [0.2839187315269263, 0.9588483477005846],
 [0.2848216774098792, 0.9585805193500557],
 [0.28490737963640583, 0.9585550505989299],
 [0.2828468469575132, 0.9591650854603669],
 [0.27424240625163093, 0.9616605963713577],
 [0.2633256370467511, 0.9647070067507143],
 [0.2858699884915783, 0.958268412126699],
 [0.2855359490356974, 0.9583679991570482],
 [0.2869073234121886, 0.9579583434432073],
 [0.28713457785972973, 0.9578902516454141],
 [0.29024666225192924, 0.9569518666325985],
 [0.29034787420366054, 0.9569211628684023],
 [0.2906651705296417, 0.9568248317434985],
 [0.28800838692634273, 0.9576278865300896],
 [0.2916117358986437, 0.9565367716330511],
 [0.29116161578269606, 0.9566738804288585],
 [0.29212614799498027, 0.956379795718006],
 [0.2920730352508251, 0.9563960173899567],
 [0.2931899571800256, 0.9560542081957355],
 [0.29565517699940536, 0.9552947274603008],
 [0.2942486550630724, 0.9557288993190344],
 [0.2956837280711308, 0.9552858906913456],
 [0.296185537019156, 0.9551304244241591],
 [0.2961898556879876, 0.9551290851960949],
 [0.295711663547997, 0.9552772435485294],
 [0.29631921181100485, 0.9550889616740971],
 [0.2962812500611317, 0.9551007385936905],
 [0.2940858488375231, 0.9557790087219501],
 [0.2922552926885215, 0.9563403389461026],
 [0.297455352885983 2, 0.9547357294243655],
 [0.2975019851605442, 0.9547211995266133],
 [0.297928771152136, 0.9545881034874562],
```
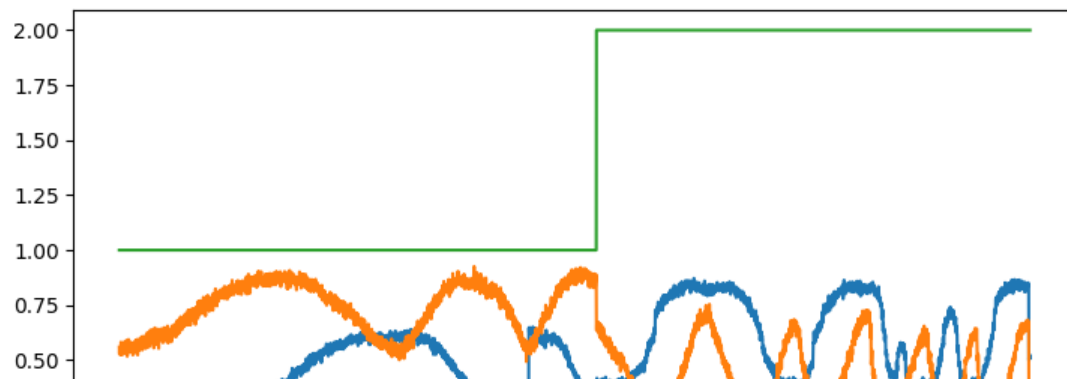
```
       [0.2982749931359468, 0.9544799780350297],
       [0.2971421220725826, 0.9548332625595961],
       [0.29847650945703846, 0.9544169808329808],
       [0.2984726016411747, 0.9544182029223608],
       [0.296753085796801, 0.9549542429195634],
       [0.2916566072042332, 0.9565230909257316],
       [0.28633736098090873, 0.95812886174381],
       [0.2995681609197765, 0.9540749011283878],
       [0.2998554748750533, 0.9539846404358177],
```
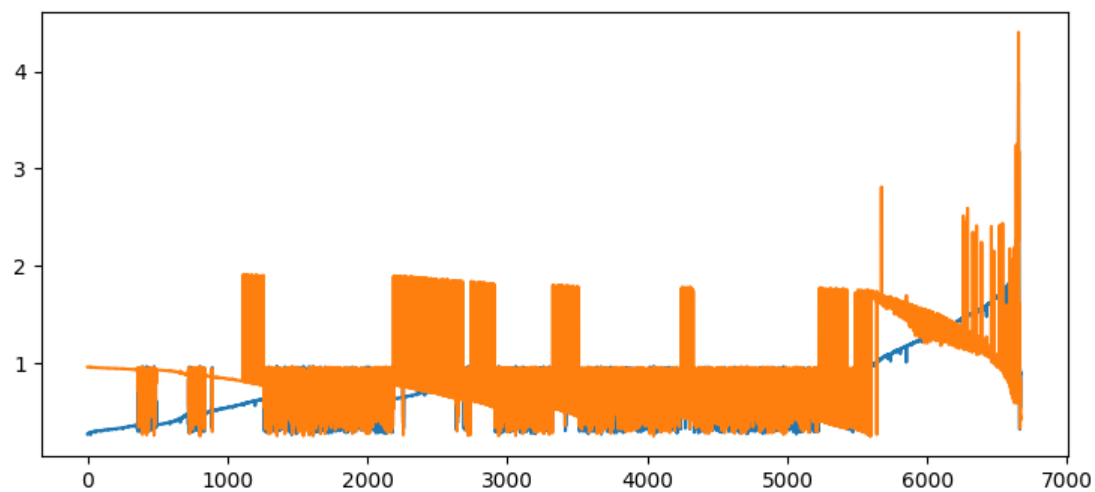
```
 1
 2 import numpy as np
 3 import matplotlib.pyplot as plt
 4
 5 plt.rcParams["figure.figsize"] = [7.50, 3.50]
 6 plt.rcParams["figure.autolayout"] = True
 7
 8 x = results[0]
 9 y = results[1]
10
11 plt.title("Line graph")
12 plt.plot(x, y, color="red")
13
14 plt.show()
```
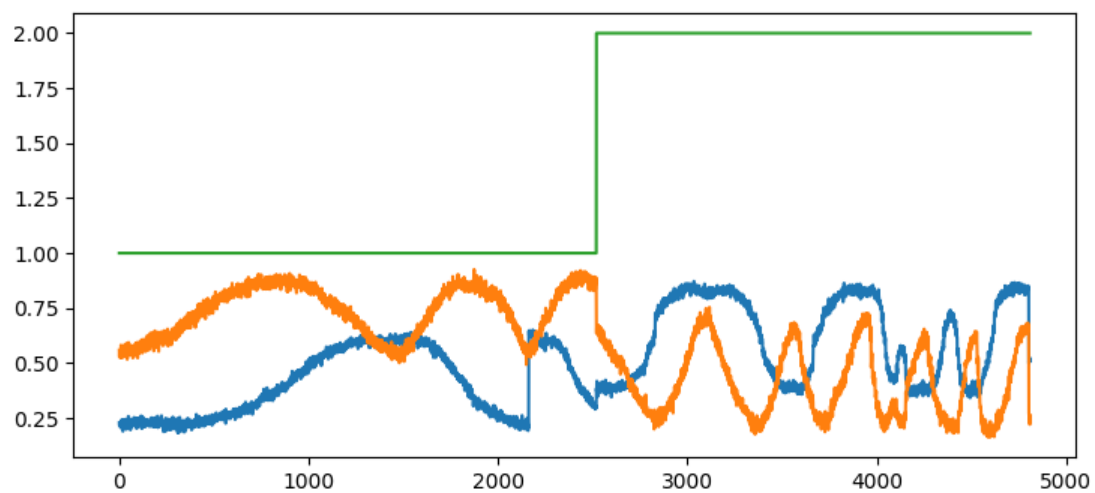


```
1 plt.plot(results)
2 plt.show()
3
```
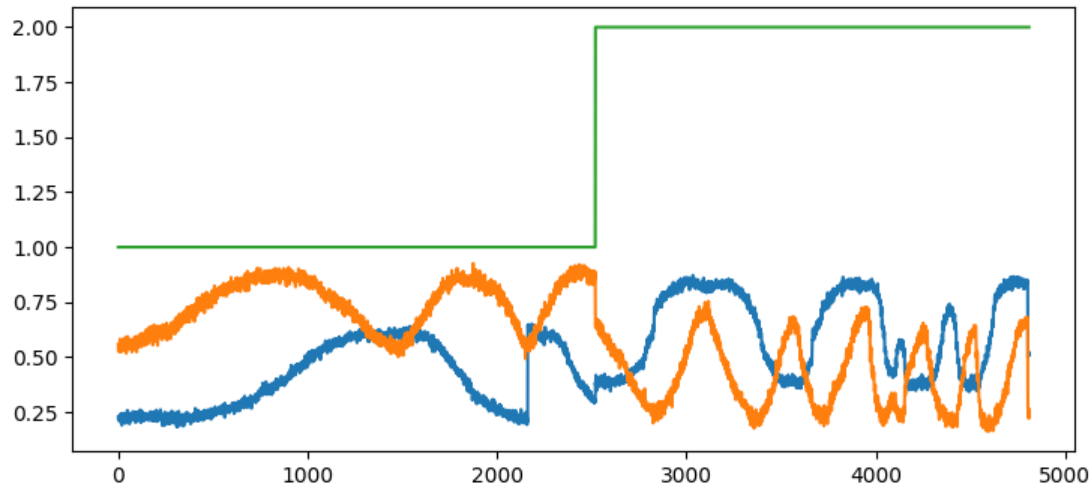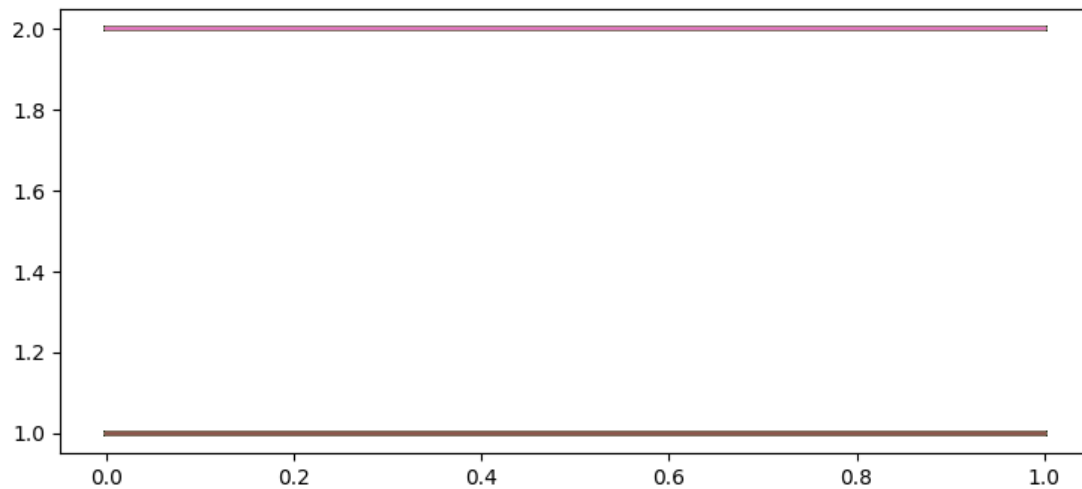
```
1 plt.plot(f_list)
2 plt.show()
```



```
1 plt.plot(data)
2 plt.show()
```
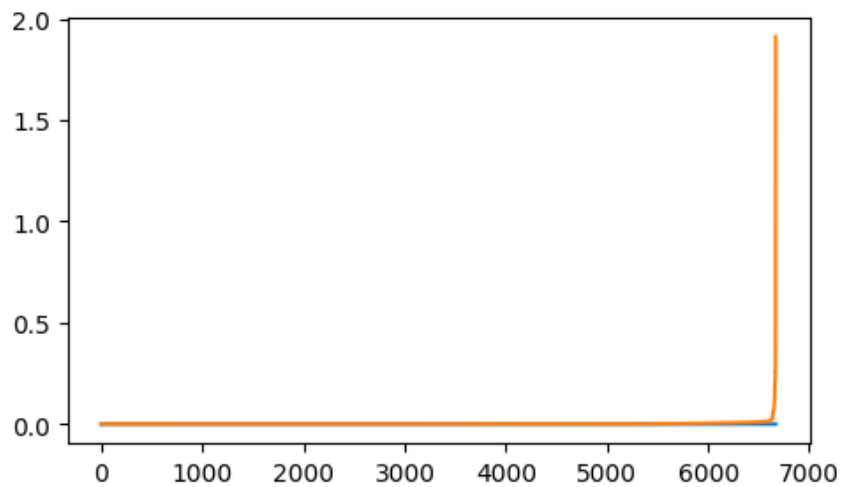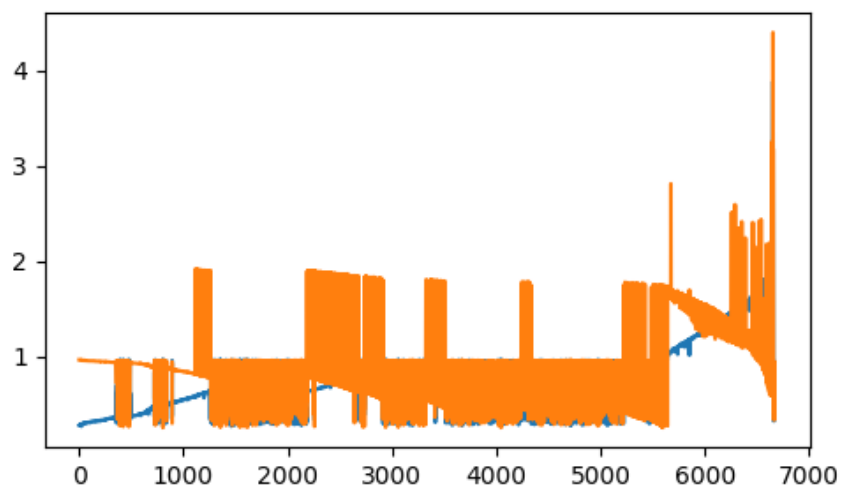
```
1 plt.plot(data1)
2 plt.show()
```



```
1 plt.plot(frames)
2 plt.show()
```



```
1 from sklearn.neighbors import NearestNeighbors
2 neighb = NearestNeighbors(n_neighbors=2)
3 nbrs=neighb.fit(f_list)
4 distances,indices=nbrs.kneighbors(f_list)
5 distances = np.sort(distances, axis = 0)
6 distances = distances[:,]
7 plt.rcParams['figure.figsize'] = (5,3)
8 plt.plot(distances)
9 plt.show()
```

```
1 df = pd.DataFrame(f_list, columns = ['x','y'])
2 plt.plot(df)
3 plt.show()
```



1