Part I: Exercises **Exercise 1** In [1]: import pandas from scipy.cluster.hierarchy import linkage from scipy.cluster.hierarchy import dendrogram

LAB 4: Clustering (Unsupervised learning)

d = [0.3, 0.4, 0.7, 0.5, 0.8, 0.45] # d should be a condensed distance matrix, like an object returned by pdist() ModuleNotFoundError Traceback (most recent call last) ~\AppData\Local\Temp/ipykernel_16080/2097474627.py in <module> 1 import pandas ---> 2 from scipy.cluster.hierarchy import linkage 3 from scipy.cluster.hierarchy import dendrogram pdist() function ModuleNotFoundError: No module named 'scipy' # a) complete linkage exo1_hc_complete = linkage(d, "complete") dendrogram(exo1_hc_complete,labels=[1,2,3,4]) 'dcoord': [[0.0, 0.3, 0.3, 0.0], [0.0, 0.45, 0.45, 0.0], [0.3, 0.8, 0.8, 0.45]],

4 d = [0.3, 0.4, 0.7, 0.5, 0.8, 0.45] # d should be a condensed distance matrix, like an object returned by 'icoord': [[5.0, 5.0, 15.0, 15.0], [25.0, 25.0, 35.0, 35.0], [10.0, 10.0, 30.0, 30.0]],

'ivl': [1, 2, 3, 4], 'leaves': [0, 1, 2, 3]}

In []: Out[]: {'color_list': ['g', 'r', 'b'], # b) complete single exo1 hc single = linkage(d, "single")

dendrogram(exo1 hc single, labels=[1,2,3,4])

In []: 'dcoord': [[0.0, 0.3, 0.3, 0.0], [0.0, 0.4, 0.4, 0.3],[0.0, 0.45, 0.45, 0.4]],'icoord': [[25.0, 25.0, 35.0, 35.0], [15.0, 15.0, 30.0, 30.0], [5.0, 5.0, 22.5, 22.5]], 'ivl': [4, 3, 1, 2], 'leaves': [3, 2, 0, 1]}

Out[]: {'color_list': ['g', 'b', 'b'], 0.4 0.3

0.2 0.1

0.0

[[0]] [0] [1] [1]] [[0]] [0] [0] [1]]

print(df)

print("JOINS")

print(df)

Exercise 2

JOINS

1

0.0

2.0

4.0

0.0

2.0

array([0.3, 0.45, 0.8])

X1 = [1, 1, 0, 5, 6, 4]X2 = [4, 3, 4, 2, 2, 0]

from random import choice

def calculate centroids(X,labels):

n1+=1

n2+=1

else:

for i in range(0,len(X[0])): if labels[i] == 1:

labels = [choice(l) for x in range(0,len(X1))]

centroid1, centroid2, n1, n2=[0,0], [0,0], 0,0

centroid1[0]+=X[0][i] centroid1[1]+=X[1][i]

centroid2[0]+=X[0][i]centroid2[1]+=X[1][i]

c1,c2 =calculate_centroids(X,labels)[0],calculate_centroids(X,labels)[1]

return (math.sqrt((x[0] - centroid[0])**2 + (x[1]-centroid[1])**2))

if (dist([X[0][i],X[1][i]], centroid1) < dist([X[0][i],X[1][i]], centroid2)):</pre>

centroid1[0]=centroid1[0]*1.0/n1 centroid1[1]=centroid1[1]*1.0/n1 centroid2[0]=centroid2[0]*1.0/n2 centroid2[1]=centroid2[1]*1.0/n2 return[centroid1,centroid2]

plt.plot(c1[0],c1[1],"*",color="black") plt.plot(c2[0],c2[1],"*",color="black")

[<matplotlib.lines.Line2D at 0x10ab1e290>]

1 = [1, 2]

print(labels)

#c) centroids

print(c1) print(c2)

plt.plot(X1, X2, "o")

[2, 1, 2, 1, 2, 2]

plt.plot(X1, X2, "o")

import matplotlib.pyplot as plt

[<matplotlib.lines.Line2D at 0x10a57e310>]

0

1

#c) (1,2), (3,4)

#d) (1, 2, 3), (4)

from scipy.cluster.hierarchy import cut tree print(cut_tree(exo1_hc_complete, n_clusters = 2))

print(cut tree(exo1 hc single, n clusters = 2))

column_labels = ["sub-cluster", "sub-cluster", "distance", "cluster size"]

0.30

0.45

0.80

0.30

0.40

0.45

exo1 hc complete[:,2] # to get the distances in the order the joins took place

2.0

2.0

4.0

3.0

4.0

df = pandas.DataFrame(exo1_hc_complete, columns=column_labels)

df = pandas.DataFrame(exo1_hc_single, columns=column_labels)

sub-cluster sub-cluster distance cluster size

1.0

3.0

5.0

#print("sub-clusters", "distance", "cluster size")

sub-cluster sub-cluster distance cluster size

1.0

4.0 5.0

In []:

In []:

In []:

In []:

Out[]:

In []:

In []:

Out[]:

In []:

Out[]:

In []:

In []:

In []:

Out[]:

In []:

4.0 3.5 3.0 2.5 2.0 1.5 1.0 0.5 0.0

4.0 3.5

2.5 2.0 1.5 1.0 0.5 0.0

import math

def dist(x,centroid):

return(labels)

X = [X1, X2]

print(labels)

[2, 2, 2, 1, 1, 1]

last_labels = []

dist([X1[0], X2[0]], centroid1)

def assign_labels(X, centroid1, centroid2):

labels.append(1)

labels.append(2)

labels = assign_labels(X, centroid1, centroid2)

centroid1 = calculate_centroids(X,labels)[0] centroid2 = calculate_centroids(X,labels)[1]

labels = assign_labels(X, centroid1, centroid2)

for i in range(len(X[0])):

while (not(last_labels == labels)):

plt.plot(X[0][i], X[1][i], "o", color=col)

[<matplotlib.lines.Line2D at 0x10a46ccd0>]

plt.plot(centroid1[0], centroid1[1], "*", color="black") plt.plot(centroid2[0],centroid2[1],"*",color="black")

last_labels = labels

print(centroid1) print(centroid2)

[5.0, 1.3333333333333333333333]

for i in range(len(X[0])): if labels[i] ==1: col="red"

col="blue"

else:

4.0 3.5 3.0 2.5 2.0 1.5 1.0 0.5 0.0