

```
"""
Created on Fri Oct 28 10:42:10 2022
@author:Brandon Botzer btb5103
```

#### Data Set Information:

The examined group comprised kernels belonging to three different varieties of wheat: Kana, Rosa and Canadian. 70 elements each, randomly selected for the experiment. High quality visualization of the internal kernel structure was detected using a soft X-ray technique. It is non-destructive and considerably cheaper than other more sophisticated imaging techniques like scanning microscopy or laser technology. The images were recorded on 13x18 cm X-ray KODAK plates. Studies were conducted using combine harvested wheat grain originating from experimental fields, explored at the Institute of Agrophysics of the Polish Academy of Sciences in Lublin.

#### Attribute Information:

To construct the data, seven geometric parameters of wheat kernels were measured:

1. area A,  $C = 4 \cdot \pi \cdot A / P^2$ ,
2. perimeter P,
3. compactness C =  $4 \cdot \pi \cdot A / P^2$ ,
4. length of kernel,
5. width of kernel,
6. asymmetry coefficient
7. length of kernel groove.

#### Relevant Papers:

M. Charytanowicz, J. Niewczas, P. Kulczycki, P.A. Kowalski, S. Lukasik, S. Zak, 'A Complete Gradient Clustering Algorithm for Features Analysis of X-ray Images', in: Information Technologies in Biomedicine, Ewa Pietka, Jacek Kawa (eds.), Springer-Verlag, Berlin-Heidelberg, 2010, pp. 15-24.

Please use this data to finish the following tasks.

1. Explore the data set (10 points)
2. Use K-means clustering to group the seed data. (30 points)
3. Use different linkage type for Hierarchical clustering to the seed data, which linkage type give the best result? (30 points)
4. Use DBSCAN clustering group the seed data and find the best epses and min\_samples value. (30 points)

```
"""
Created on Fri Oct 28 10:42:10 2022
@author:Brandon Botzer btb5103

Data Set Information:
The exam
ined group comprised kernels belonging to three different varieties \n wheat: Kana, Rosa and Canadian, 70 ele
ments each, randomly selected for \n the experiment. High quality visualization of the internal kernel structure
\n was detected using a soft X-ray technique. It is non-destructive and considerably cheaper than other more s
ophisticated imaging techniques like \n scanning microscopy or laser technology. The images were recorded on \n
13x18 cm X-ray KODAK plates. Studies were conducted using combine harvested \n wheat grain originating from exper
imental fields, explored at the Institute \n of Agrophysics of the Polish Academy of Sciences in Lublin. \n \n \n
Attribute Information:
\n\n\n To construct the data, seven geometric parameters of wheat kernels were \n measured: \n
1. area A, \n 2. perimeter P, \n 3. compactness C = 4 \cdot \pi \cdot A / P^2, \n 4. length of kernel, \n 5. w
idth of kernel, \n 6. asymmetry coefficient \n 7. length of kernel groove. \n \n \n
Relevant Papers:
\n\n M. Charytanowicz, J. Niewczas, P. Kulczycki, P.A. Kowalski, S. Lukasik, S. Zak, \n 'A Complete Gradient Clustering Algorithm for Features Analysis of
X-ray Images', in: Information Technologies in Biomedicine, Ewa Pietka, Jacek Kawa (eds.), Springer-Verlag, Berlin-Heidelberg, 2010, pp. 15-24. \n \n \n
Please use this data to finish the following tasks.
\n\n 1. Explore the data set (10 points) \n 2. Use K-means clustering to g
roup the seed data. (30 points) \n 3. Use different linkage type for Hierarchical clustering to the seed dat
a, \n which linkage type give the best result? (30 points) \n 4. Use DBSCAN clustering group the seed data
and find the best epses \n and min_samples value. (30 points) \n \n \n
"""
```

```
In [2]: import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import metrics

#0. Read the data in
#change the directory
os.chdir('readpath')
#This data set does not have a header for the column names
#set the column names for the data frame
column_names = ['area', 'perimeter', 'compactness', 'length_kernel',
                'width_kernel', 'asymmetry', 'length_groove', 'type']
#read in the csv data, set header=0 and names=column_names to
have proper naming for the column information
df = pd.read_csv('seeds_dataset.csv', header=0, names=column_names)

#In case I want to split the class into dummy variables
dummy_class = pd.get_dummies(df['type'])
#drop class so I can use the dummy class instead
df_dummy = df.drop('type', axis = 1)
#right Join the dummy classes to the seed dataframe
df_dummy = df_dummy.join(dummy_class, how = 'right')

#bring the seed variables by themselves as a data frame
seeds = df.iloc[:, :7]
```

```
In [3]: #1. Explore the data set (10 points)
corr_val = seeds.corr()

#plot a correlation matrix
plt.matshow(corr_val)
plt.title('Correlation Matrix')
plt.xticks(range(7), list(seeds.columns))
plt.yticks(range(7), list(seeds.columns))

#get some info about the seeds data frame
print("\nData frame information:\n")
print(seeds.info())
print()
print(seeds.describe())

#the data should be clean but I'll still check for duplicates
dups = seeds.duplicated()
print("\nAre there any duplicates?\n" + str(dups.max()))

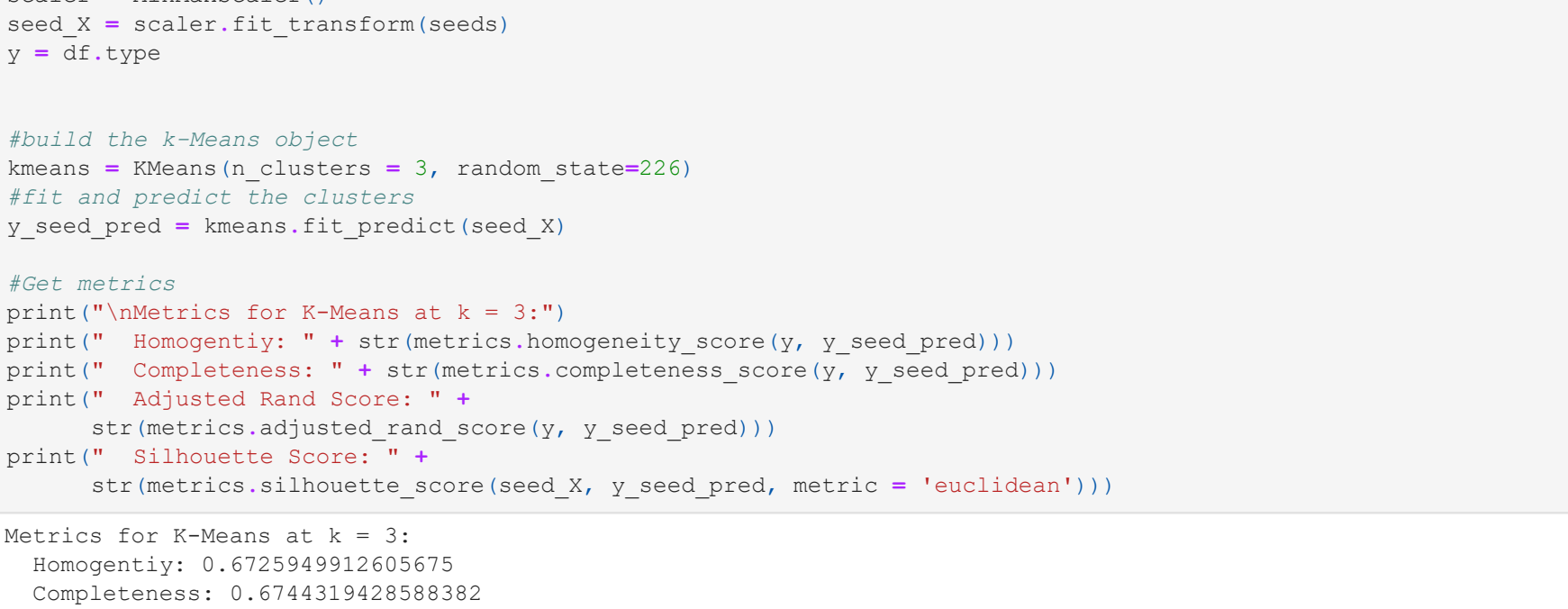
Data frame information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 209 entries, 0 to 208
Data columns (total 7 columns):
#   column      Non-Null Count  Dtype
---  ---
0   area         209 non-null             float64
1   perimeter    209 non-null             float64
2   compactness  209 non-null             float64
3   length_kernel 209 non-null             float64
4   width_kernel 209 non-null             float64
5   asymmetry    209 non-null             float64
6   length_groove 209 non-null             float64
dtypes: float64(7)
memory usage: 11.6 KB
None

          area  perimeter  compactness  length_kernel  width_kernel  \
count  209.000000  209.000000  209.000000  209.000000  209.000000
mean     14.845550   14.557943   0.870999    0.562789    0.258349
std       2.816545    1.308949   0.023686    0.444029    0.378603
min       0.159000    0.410000    0.808100    4.899000    2.630000
25%      12.260000   13.450000   0.856700    5.262000    2.941000
50%      14.340000   14.290000   0.873500    5.520000    3.232000
75%      17.320000   15.730000   0.887900    5.980000    3.562000
max      21.180000   17.250000   0.918300    6.675000    4.033000

          asymmetry  length_groove
count  209.000000  209.000000
mean     3.707278    5.408971
std     1.503457    0.492487
min     0.765300    4.519300
25%     2.587000    5.045000
50%     3.600000    5.224000
75%     4.773000    5.877000
max     8.456000    6.550000
```

Are there any duplicates?

False



In [4]: #2. Use K-means clustering to group the seed data. (30 points)

```
#I will start this with k = 3
#I will later show using the elbow method that k = 3 is correct

from sklearn.cluster import KMeans
from sklearn.preprocessing import MinMaxScaler

#Normalize independent variables for the seed data

#make the scaler object
scaler = MinMaxScaler()
seed_X = scaler.fit_transform(seeds)
y = df.type

#build the k-Means object
kmeans = KMeans(n_clusters = 3, random_state=226)
#fit and predict the clusters
y_pred = kmeans.fit_predict(seed_X)

#Get metrics
print("\nMetrics for K-Means at k = 3:")
print(" Homogeneity: " + str(metrics.homogeneity_score(y, y_pred)))
print(" Completeness: " + str(metrics.completeness_score(y, y_pred)))
print(" Adjusted Rand Score: " + str(metrics.adjusted_rand_score(y, y_pred)))
print(" Silhouette Score: " + str(metrics.silhouette_score(seed_X, y_pred, metric = 'euclidean')))
```

Metrics for K-Means at k = 3:

Homogeneity: 0.6725949912605675

Completeness: 0.674419428508392

Adjusted Rand Score: 0.442461490237547

Silhouette Score: 0.42119810498616345

```
In [5]: #try to create a plot that shows the clustering worked
#try this for all combinations of the variables to see good
#separating distinguishers

#I'll comment the loop out later and only show one plot of this
"""
for i in range(0,7):
    for j in range(0,7):
        if i != j:
            centers = kmeans.cluster_centers_
            centers_a = centers[:, i]
            centers_b = centers[:, j]

            plt.figure()
            # predicted categories (c = y_pred).
            plt.scatter(seed_X[:, i], seed_X[:, j], c = y_pred,
                        label = 'Predicted')
            # Plot the cluster centroids with 'x' and color is red
            plt.scatter(centers_a, centers_b, s=100, c = 'r',
                        marker = 'x', label = 'Cluster Centers')
            plt.title('Clustering results')
            plt.xlabel(seeds.columns[i])
            plt.ylabel(seeds.columns[j])
            plt.legend()
            """
```

Out[5]:

```
"""for i in range(0,7):
    for j in range(0,7):
        if i != j:
            centers = kmeans.cluster_centers_
            centers_a = centers[:, i]
            centers_b = centers[:, j]

            plt.figure()
            # predicted categories (c = y_pred).
            plt.scatter(seed_X[:, i], seed_X[:, j], c = y_pred,
                        label = 'Predicted')
            # Plot the cluster centroids with 'x' and color is red
            plt.scatter(centers_a, centers_b, s=100, c = 'r',
                        marker = 'x', label = 'Cluster Centers')
            plt.title('Clustering results')
            plt.xlabel(seeds.columns[i])
            plt.ylabel(seeds.columns[j])
            plt.legend()
            """
```

In [6]: #This is a good plot to show the clustering

```
"""for i in range(0,7):
    for j in range(0,7):
        if i != j:
            centers = kmeans.cluster_centers_
            centers_a = centers[:, i]
            centers_b = centers[:, j]

            plt.figure()
            # predicted categories (c = y_pred).
            plt.scatter(seed_X[:, i], seed_X[:, j], c = y_pred,
                        label = 'Predicted')
            # Plot the cluster centroids with 'x' and color is red
            plt.scatter(centers_a, centers_b, s=100, c = 'r',
                        marker = 'x', label = 'Cluster Centers')
            plt.title('Clustering results')
            plt.xlabel(seeds.columns[i])
            plt.ylabel(seeds.columns[j])
            plt.legend()
            """
```

Out[6]:

<matplotlib.legend.Legend at 0x11781f64c>



In [7]: #Use the elbow method to show that k = 3 is the proper choice

```
#computing distances from two inputs
from scipy.spatial.distance import cdist

#set up a range of k values
k_vals = range(1,10)

#set an empty list for the mean dispersions
meanDisp = []

for k in k_vals:
    #build the K-Means object with the variable n_clusters
    kmeans = KMeans(n_clusters=k)
    #fit the data with K-means
    kmeans.fit(seed_X)

    #from the sklearn book (Hacking 2nd ed. pg. 208)
    #find the average minimum distance to the cluster centers
    #shows if you've done a good job with clustering
    #if K = number of points then this would be zero
    meanDisp.append(sum(np.min(cdist(seed_X, kmeans.cluster_centers_,
                                    axis = 1)) / seed_X.shape[0]))

#plot the elbow plot
plt.figure()
plt.plot(k_vals, meanDisp, 'bx-')
plt.xlabel('K')
plt.ylabel('Average Dispersion')
plt.title('Elbow plot to show best K value')
```

C:\Users\btb51\Anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1036: UserWarning: KMeans is known to have

memory leak issues. Please use MiniBatchKMeans, where there are less chunks than available threads. You can avoid it by set

ting the environment variable OMP\_NUM\_THREADS=1.

warnings.warn('KMeans is known to have memory leak issues. Please use MiniBatchKMeans, where there are less chunks than available threads. You can avoid it by set

ting the environment variable OMP\_NUM\_THREADS=1.', UserWarning)

Out[7]:

Text(0.5, 1.0, 'Elbow plot to show best K value')



In [8]: #3. Use different linkage type for Hierarchical clustering to the seed data,

# which linkage type give the best result? (30 points)

```
from sklearn.cluster import AgglomerativeClustering

#Using linkage = 'average'
hier1_seed_pred = AgglomerativeClustering(n_clusters=3,
                                          affinity='euclidean',
                                          linkage='average').fit_predict(seed_X)

#Get metrics
print("\nMetrics for Hierarchical using Average Linkage:")
print(" Homogeneity: " + str(metrics.homogeneity_score(y, hier1_seed_pred)))
print(" Completeness: " + str(metrics.completeness_score(y, hier1_seed_pred)))
print(" Adjusted Rand Score: " + str(metrics.adjusted_rand_score(y, hier1_seed_pred)))
print(" Silhouette Score: " + str(metrics.silhouette_score(seed_X, hier1_seed_pred, metric = 'euclidean')))
```

Metrics for Hierarchical using Average Linkage:

Homogeneity: 0.49095617606258773

Completeness: 0.645373983245087

Adjusted Rand Score: 0.5079710294939454

Silhouette Score: 0.33007475633536437

In [9]: #Using linkage = 'complete'

```
hier2_seed_pred = AgglomerativeClustering(n_clusters=3,
                                          affinity='euclidean',
                                          linkage='complete').fit_predict(seed_X)

#Get metrics
print("\nMetrics for Hierarchical using Complete Linkage:")
print(" Homogeneity: " + str(metrics.homogeneity_score(y, hier2_seed_pred)))
print(" Completeness: " + str(metrics.completeness_score(y, hier2_seed_pred)))
print(" Adjusted Rand Score: " + str(metrics.adjusted_rand_score(y, hier2_seed_pred)))
print(" Silhouette Score: " + str(metrics.silhouette_score(seed_X, hier2_seed_pred, metric = 'euclidean')))
```

Metrics for Hierarchical using Complete Linkage:

Homogeneity: 0.608607015344328

Completeness: 0.6345373983245087

Adjusted Rand Score: 0.5651989293798214

Silhouette Score: 0.362782119361832

In [10]: #Using linkage = 'ward'

```
hier3_seed_pred = AgglomerativeClustering(n_clusters=3,
                                          affinity='euclidean',
                                          linkage='ward').fit_predict(seed_X)

#Get metrics
print("\nMetrics for Hierarchical using Ward Linkage:")
print(" Homogeneity: " + str(metrics.homogeneity_score(y, hier3_seed_pred)))
print(" Completeness: " + str(metrics.completeness_score(y, hier3_seed_pred)))
print(" Adjusted Rand Score: " + str(metrics.adjusted_rand_score(y, hier3_seed_pred)))
print(" Silhouette Score: " + str(metrics.silhouette_score(seed_X, hier3_seed_pred, metric = 'euclidean')))
```

Metrics for Hierarchical using Ward Linkage:

Homogeneity: 0.4651550494911663

Completeness: 0.4591566372116655

Adjusted Rand Score: 0.6521911239329006

Silhouette Score: 0.3906386643092875

In [11]: #Based on the metrics run, the Ward Linkage performs the best for clustering

#this data set having the highest silhouette score (among other values)

```
#Plot the different Hierarchical Clusterings linkage types
plt.figure(figsize = (12,12))
plt.subplot(221)
plt.scatter(seed_X[:, 6], seed_X[:, 4], c = y)

plt.title('Clustering results')
plt.xlabel('Width_Kernel')
plt.ylabel('Length_Groove')

plt.subplot(222)
plt.scatter(seed_X[:, 6], seed_X[:, 4], c = hier1_seed_pred)
plt.title('Average Linkage')
plt.xlabel('Width_Kernel')
plt.ylabel('Length_Groove')

plt.subplot(223)
plt.scatter(seed_X[:, 6], seed_X[:, 4], c = hier2_seed_pred)
plt.title('Complete Linkage')
plt.xlabel('Width_Kernel')
plt.ylabel('Length_Groove')

plt.subplot(224)
plt.scatter(seed_X[:, 6], seed_X[:, 4], c = hier3_seed_pred)
plt.title('Ward Linkage')
plt.xlabel('Width_Kernel')
plt.ylabel('Length_Groove')

plt.subplots_adjust(wspace = 0.2)
```

Out[11]:

Clustering results

Average Linkage

Complete Linkage

Ward Linkage

In [12]: #4. Use DBSCAN clustering group the seed data and find the best epses

# and min\_samples value. (30 points)

```
from sklearn.cluster import DBSCAN

#Build the model and fit predict the clusters
db_y_seed_pred = DBSCAN(eps=0.2, min_samples=5).fit_predict(seed_X)
```

#Get metrics

print("\nMetrics for DBSCAN with eps = 0.2 and min\_samples = 5:")

print(" Homogeneity: " + str(metrics.homogeneity\_score(y, db\_y\_seed\_pred)))

print(" Completeness: " + str(metrics.completeness\_score(y, db\_y\_seed\_pred)))

print(" Adjusted Rand Score: " + str(metrics.adjusted\_rand\_score(y, db\_y\_seed\_pred)))

print(" Silhouette Score: " + str(metrics.silhouette\_score(seed\_X, db\_y\_seed\_pred, metric = 'euclidean')))

Metrics for DBSCAN with eps = 0.2 and min\_samples = 5:

Homogeneity: 0.4651550494911663

Completeness: 0.4591566372116655

Adjusted Rand Score: 0.3854448634577096

Silhouette Score: 0.08210506434215879

In [13]: #Run a search for the best 'epses' and 'min\_samples' for the DBSCAN

```
print("\n\nSearching for the best parameters for DBSCAN...\n\n")
#make an empty list to store results
results = []

#make the 'epses' and 'min_samples' grid
epses = [0.1, 0.2, 0.3, 0.4, 0.5, 1.0]
min_samples = [2, 3, 10, 15]
```

#run over the search grid manually

for n in min\_samples:

for e in epses:

#set up the DBSCAN model

model = DBSCAN(eps = e, min\_samples = n)

#Get the predicted clusters

db\_y\_seed\_pred\_temp = model.fit\_predict(seed\_X)

#find the number of clusters

n\_clusters = np.unique(model.labels\_).size

#Get the silhouette metrics

#Check if there is just one cluster, this will error the Silhouette

#so if there is only one cluster this is a bad score and will

equal 0 (assuming there is more than one cluster...)

if n\_clusters == 1:

results.append((e,n,np.nan, n\_clusters))

else:

score = metrics.silhouette\_score(seed\_X, db\_y\_seed\_pred\_temp,

metric = 'euclidean')

results.append((e, n, score, n\_clusters))

#put results in to a data frame

res\_hold = pd.DataFrame(results, columns = ['epses',

'min\_samples',

'sil\_score',

'n\_clusters'])

print("\n\nThe Silhouette scores from the first DBSCAN search:\n\n")

print(res\_hold)

#Get the 'epses', 'min\_samples' and 'sil\_score' as reference values for the loop

r\_epses = res\_hold.epses[res\_hold.sil\_score.idxmax()]

r\_min\_samp = res\_hold.min\_samples[res\_hold.sil\_score.idxmax()]

r\_sil\_score = res\_hold.sil\_score.max()

Searching for the best parameters for DBSCAN...

The Silhouette scores from the first DBSCAN search:

epses	min_samples	sil_score	n_clusters
0	0.1	0.06771	27
1	0.1	-0.097755	2
2	0.1	NaN	1
3	0.1	NaN	1
4	0.1	NaN	1
5	0.2	0.118907	5
6	0.2	0.125795	4
7	0.2	-0.193016	4
8	0.2	-0.052258	3
9	0.3	0.055722	2
10	0.3	0.055722	2
11	0.3	0.199218	2
12	0.4	NaN	1
13	0.4	NaN	1
14	0.4	NaN	1
15	0.4	0.080234	2
16	0.5	NaN	1
17	0.5	NaN	1
18	0.5	NaN	1
19	0.5	NaN	1
20	1.0	NaN	1
21	1.0	NaN	1
22	1.0	NaN	1
23	1.0	NaN	1
24	1.0	NaN	1

In [14]: #We've found the epses and min\_samples around 0.3 and 15 to be best

```
#I'll refine search to obtain better sil score
#Do this in the while loop
#Find epses and min_num from trial 1 above, then do + / - around the params
#while loop goes until sil_score results don't increase (could be local max)
# I have not accounted for local maximums (could provide inertia for this)
```

check = False

min\_samp\_inc = 1

while check == False:

#make and empty list to store results

results = []

#make the 'epses' and 'min\_samples' grid

epses\_hold = [r\_epses + (2\*epses\_inc),

r\_epses-epses\_inc,

r\_epses\*epses\_inc,

r\_epses\*(2\*epses\_inc)]

min\_samps\_hold = [r\_min\_samp + (2\*min\_samp\_inc),

min\_samps\_hold - (2\*min\_samp\_inc),

r\_min\_samp,

r\_min\_samp\*(2\*min\_samp\_inc)]

#run over the search grid manually

for e in epses\_hold:

for n in min\_samps\_hold:

#set up the DBSCAN model

model = DBSCAN(eps = e, min\_samples = n)

#Get the predicted clusters

db\_y\_seed\_pred\_temp = model.fit\_predict(seed\_X)

#find the number of clusters

n\_clusters = np.unique(model.labels\_).size

#Get the silhouette metrics

#Check if there is just one cluster, this will error the Silhouette

#so if there is only one cluster this is a bad score and will

equal 0 (assuming there is more than one cluster...)

if n\_clusters == 1:

results.append((e,n,np.nan, n\_clusters))

else:

score = metrics.silhouette\_score(seed\_X, db\_y\_seed\_pred\_temp,

metric = 'euclidean')

results.append((e, n, score, n\_clusters))

#put results in to a temp data frame

res\_temp = pd.DataFrame(results, columns = ['epses',

'min\_samples',

'sil\_score',

'n\_clusters'])

#check the reference sil score for changes

if r\_sil\_score\_temp > r\_sil\_score:

#a better score has been found. Reassign the reference vals and go again

r\_sil\_score = res\_temp.sil\_score\_temp

r\_epses = r\_epses\_temp

r\_min\_samp = r\_min\_samp\_temp

elif r\_sil\_score\_temp <= r\_sil\_score:

#no improvements have been made. Hit the flag and break the loop

check = True

best\_clusters = res\_n\_clusters[res.sil\_score.idxmax()]

print("\n\nThe Silhouette scores from the last DBSCAN search:\n\n")

print(res)

print("\n\nThe best DBSCAN values are:")

print(" epses: " + str(r\_epses))

print(" min\_samps: " + str(r\_min\_samp))

print(" Silhouette Score: " + str(r\_sil\_score))

print(" N Clusters: " + str(best\_clusters))

The Silhouette scores from the last DBSCAN search:

epses	min_samples	sil_score	n_clusters
0	0.15	9	0.06771
1	0.15	10	-0.097755
2	0.15	11	NaN
3	0.15	12	NaN
4	0.15	13	NaN
5	0.20	9	0.118907
6	0.20	10	0.125795
7	0.20	12	0.082605
8	0.20	12	0.082605
9	0.20	13	0.055205
10	0.25	11	0.346873
11	0.25	12	0.33391
12	0.25	13	0.329685
13	0.30	9	0.055722
14	0.30	10	0.055722
15	0.30	13	0.115260
16	0.35	9	0.104581
17	0.35	10	0.104581
18	0.35	12	0.104581
19	0.35	13	0.104581

The best DBSCAN values are:

epses: 0.25

min\_samples: 11

Silhouette Score: 0.3468734616603017

sil\_clusters: 3