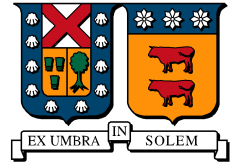




INF391 - Reconocimiento de Patrones en Minería de Datos



Tarea 1: Técnicas de *Clustering*

Francisca Ramírez

Juan Pablo Muñoz

17 de abril del 2019

Introducción

En esta tarea se exploran distintas técnicas de reconocimiento de patrones basadas en *clustering* vistas en cátedra. Para ello, se cuenta con tres pequeños *datasets* con distintas características, que servirán para contrastar la aptitud que cada técnica posee para cada caso.

Luego de la experimentación, se responden las dos preguntas conceptuales planteadas en el enunciado.

Parte I

Primero, se prepara la ingesta de datos.

In [1]:

```
1 import os.path
2 import numpy as np
3
4 def ingest_dataset(txt_dir):
5     dataset = list()
6     if os.path.exists(txt_dir):
7         with open(txt_dir, 'r') as f:
8             for line in f.readlines():
9                 data_point = line.split()
10                 x_coord, y_coord = float(data_point[0]), float(data_point[1])
11                 dataset.append([x_coord, y_coord])
12     return np.array(dataset)
```

Y se instancian los tres datasets.

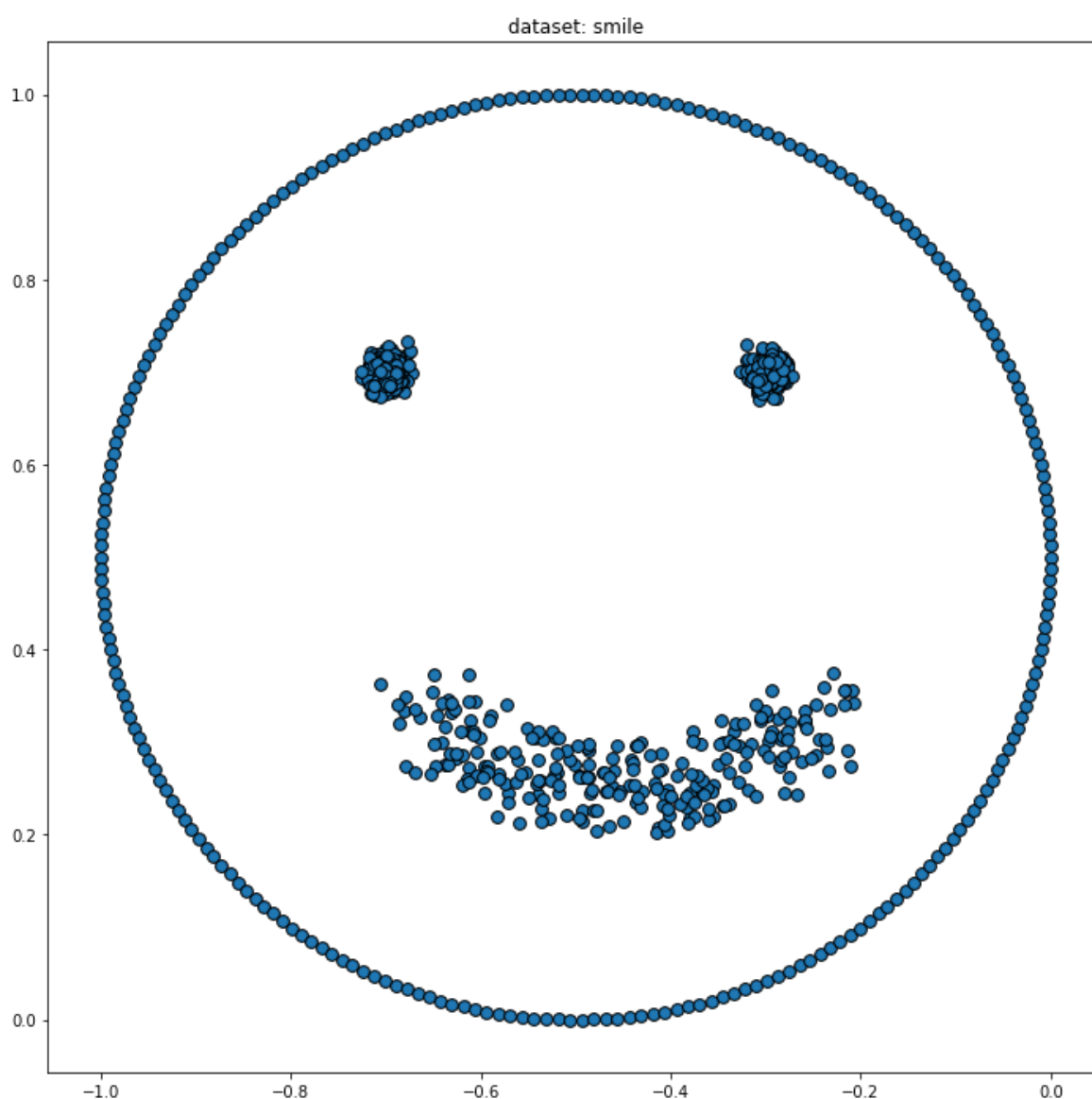
In [2]:

```
1 smile = ingest_dataset('smile.txt')
2 mouse = ingest_dataset('mouse.txt')
3 spiral = ingest_dataset('spiral.txt')
```

El hecho de que estos datasets sean 2-dimensionales nos permite visualizarlos fácilmente y obtener una clara idea cómo se distribuyen.

In [161]:

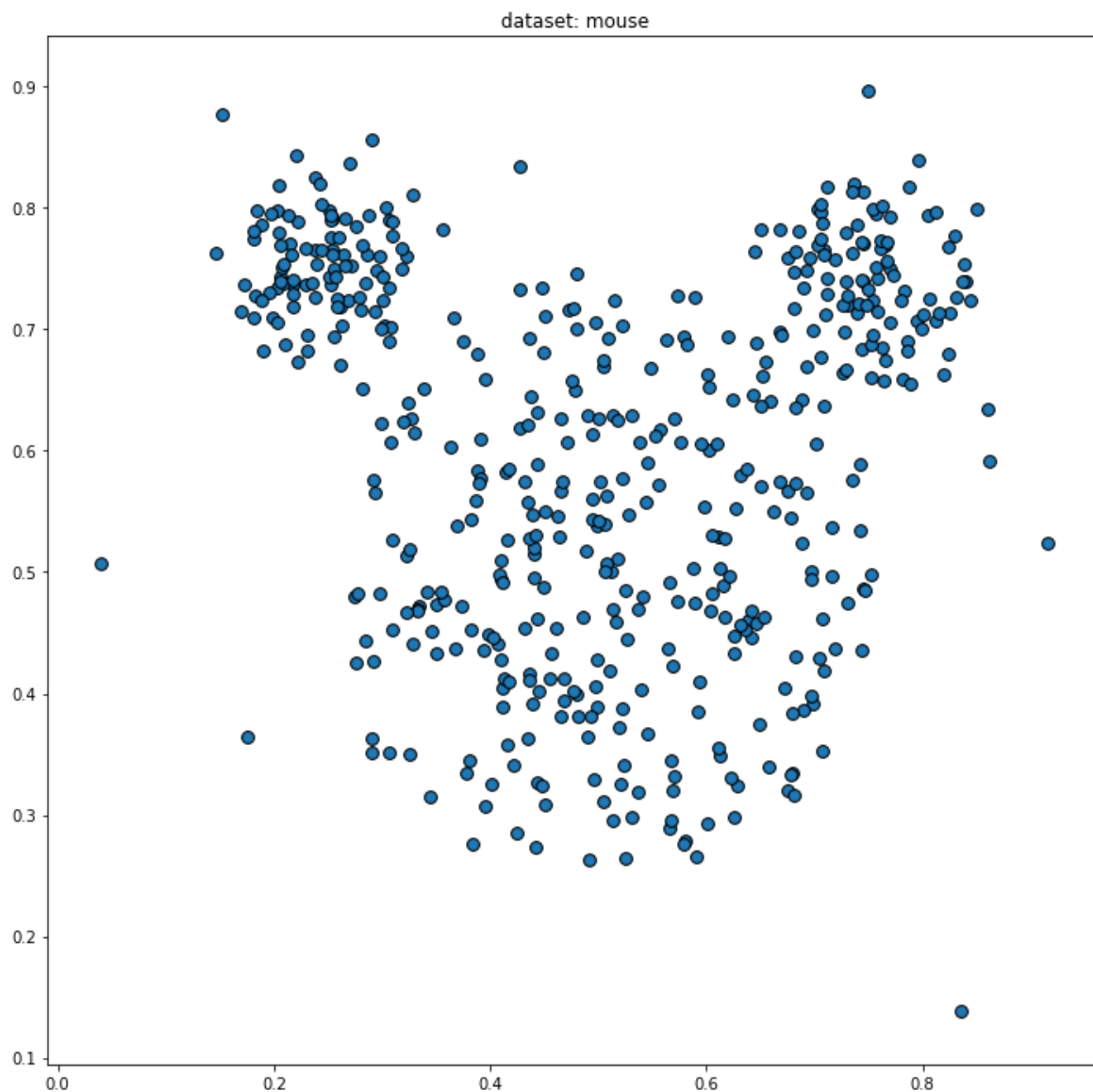
```
1 # Smile
2 plt.figure(figsize=(12,12))
3 plt.scatter(smile[:, 0], smile[:, 1], marker='o',
4             edgecolors='k', s=60)
5 plt.title('dataset: smile')
6 plt.show()
```



Del gráfico, es evidente que este dataset posee 4 clusters: dos de ellos tienen una alta densidad y tienen forma convexa o globular (los ojos), uno tiene densidad media y es semi-convexo (la sonrisa) y el último es una distribución aparentemente uniforme en forma de circunferencia, que encierra el espacio que los tres clusters anteriores ocupan.

In [154]:

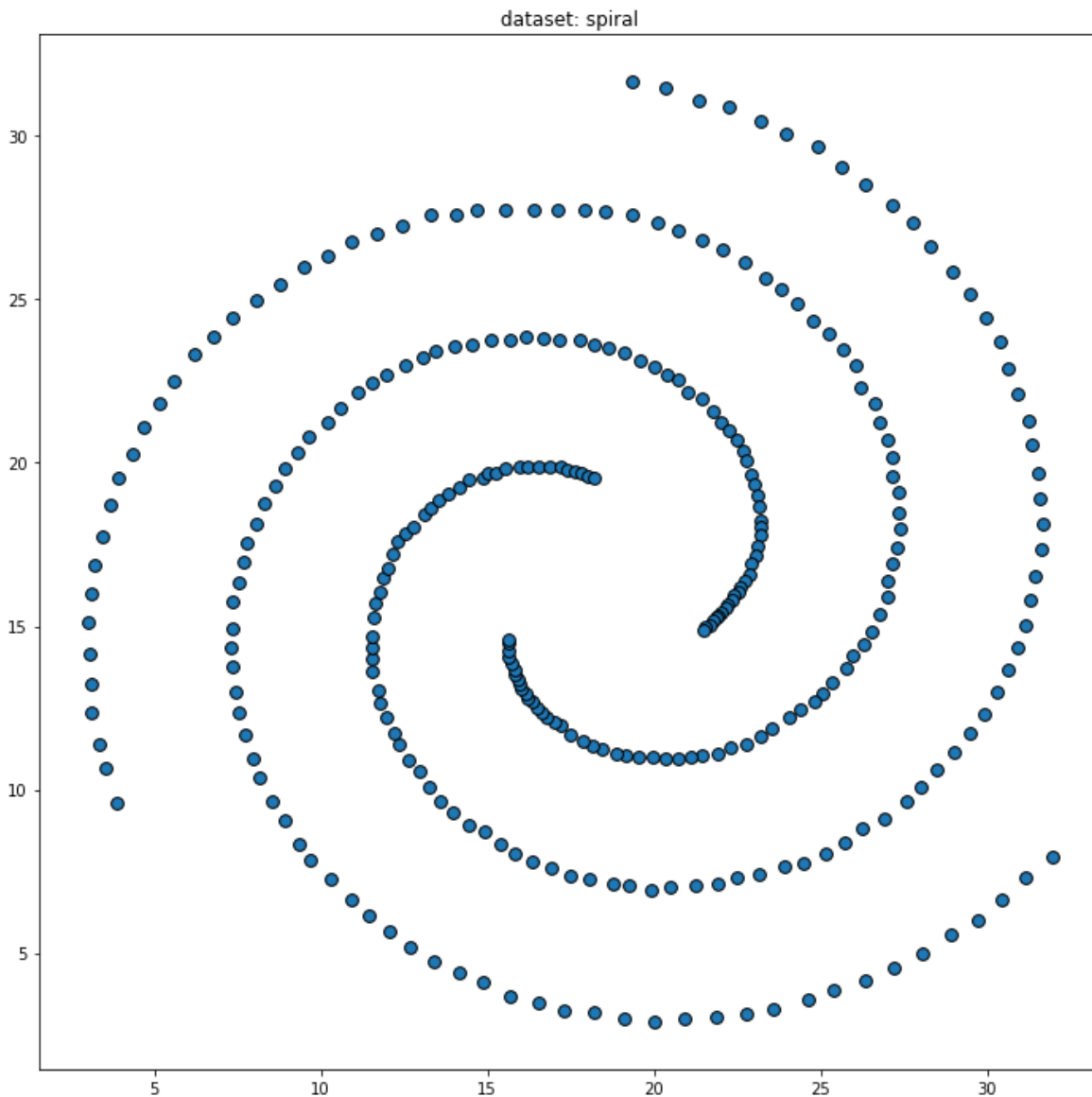
```
1 # Mouse
2 plt.figure(figsize=(12,12))
3 plt.scatter(mouse[:, 0], mouse[:, 1], marker='o',
4             edgecolors='k', s=60)
5 plt.title('dataset: mouse')
6 plt.show()
```



Este dataset posee tres clusters circulares, uno de densidad mediana (cabeza) y dos más pequeños de densidad media-alta (orejas). Los tres clusters tienen forma convexa y tienen poca a nula intersección entre sí.

In [155]:

```
1 # Spiral
2 plt.figure(figsize=(12,12))
3 plt.scatter(spiral[:, 0], spiral[:, 1], marker='o',
4             edgecolors='k', s=60)
5 plt.title('dataset: spiral')
6 plt.show()
```



La distribución de este dataset es no-convexa y resulta intuitivo distinguir cada espiral como un posible cluster distinto. En este caso, no es la densidad de los puntos lo que nos indica la presencia de un cluster, sino el patrón que estos describen.

A continuación, se procede a aplicar las técnicas de *clustering*.

(Info teórica sobre los algoritmos de clustering excepto Fuzzy se puede hallar acá: <https://scikit-learn.org/stable/modules/clustering.html> (<https://scikit-learn.org/stable/modules/clustering.html>))

1. K-Means

In [92]:

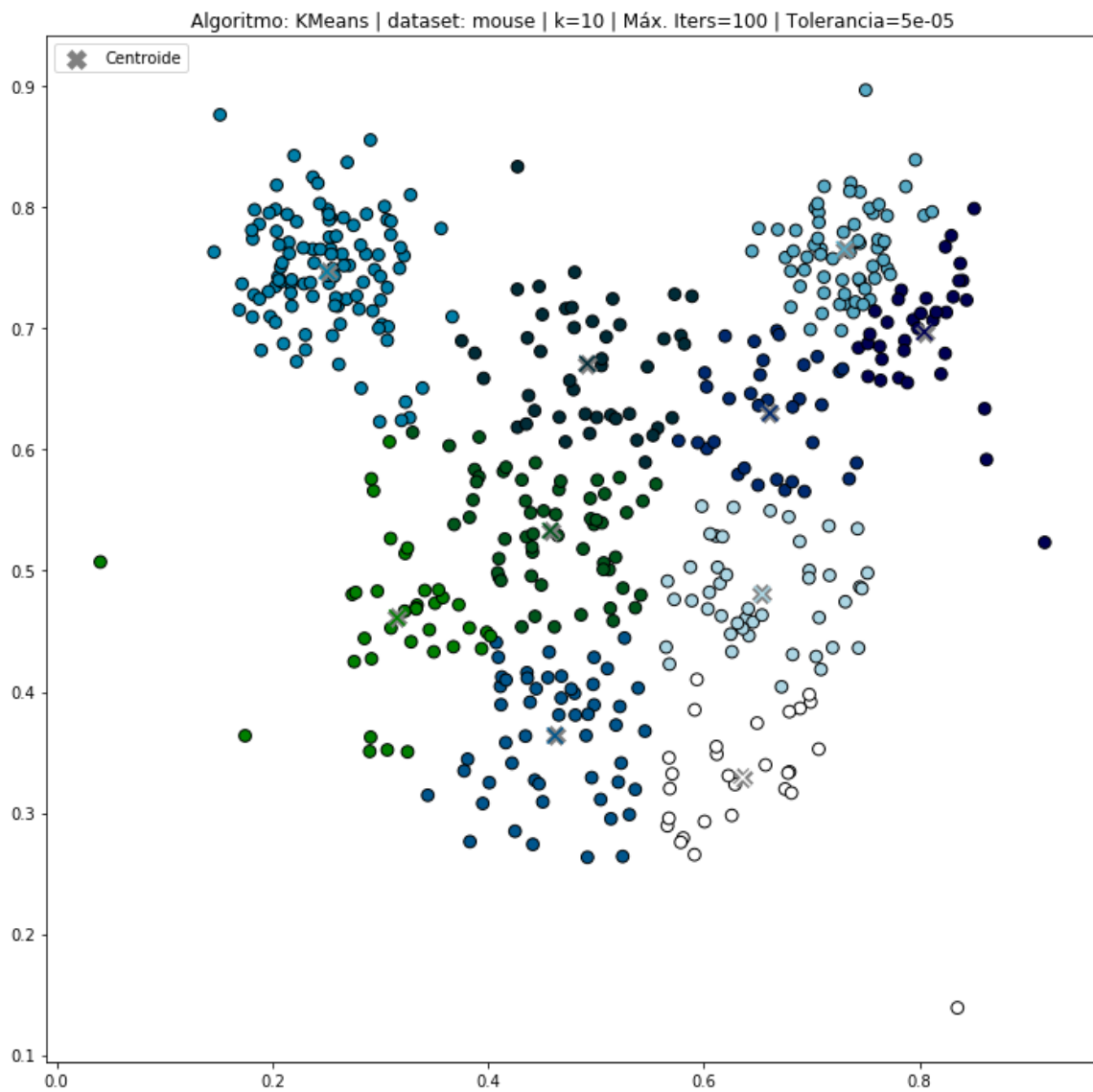
```
1 from sklearn.cluster import KMeans
2 import matplotlib.pyplot as plt
3 from ipywidgets import interact
4 from ipywidgets import FloatSlider
5
6 def apply_kmeans(dataset, k, max_iterations=300, tolerance=1e-4):
7     kmeans = KMeans(
8         n_clusters=k,
9         init='random',
10        n_init=1,
11        max_iter=max_iterations,
12        tol=tolerance,
13        random_state=0,
14    )
15    kmeans.fit(dataset)
16    return kmeans.cluster_centers_, kmeans.labels_
17
18 @interact(
19     dataset_name=['smile', 'mouse', 'spiral'],
20     k=(2,10, 1),
21     max_iterations=(10, 100, 10),
22     tolerance=FloatSlider(
23         min=5e-5,
24         max=5e-4,
25         step=5e-5,
26         continuous_update=False,
27         readout=True,
28         readout_format='.5f'
29     ),
30 )
31 def plot_kmeans(dataset_name, k, max_iterations, tolerance):
32     if dataset_name == 'smile':
33         dataset = smile
34     elif dataset_name == 'mouse':
35         dataset = mouse
36     elif dataset_name == 'spiral':
37         dataset = spiral
38     centroids, labels = apply_kmeans(dataset, k, max_iterations, tolerance)
39     plt.figure(figsize=(12,12))
40     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
41                 edgecolors='k', s=60, cmap=plt.cm.ocean)
42     plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=150,
43                 linewidths=.5, c='gray', cmap=plt.cm.ocean, label='Centroide')
44     plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=100,
45                 linewidths=2, c=list(range(len(centroids))),
46                 cmap=plt.cm.ocean)
47     plt.title('Algoritmo: KMeans | dataset: {} | k={} | Máx. Iters={} | Tolerancia={}').format(
48         dataset_name, k, max_iterations, tolerance)
```

dataset_na...

k

max_iterati...

tolerance



Análisis K-Means

Bla...

2. Agglomerative Hierarchical Clustering

In [42]:

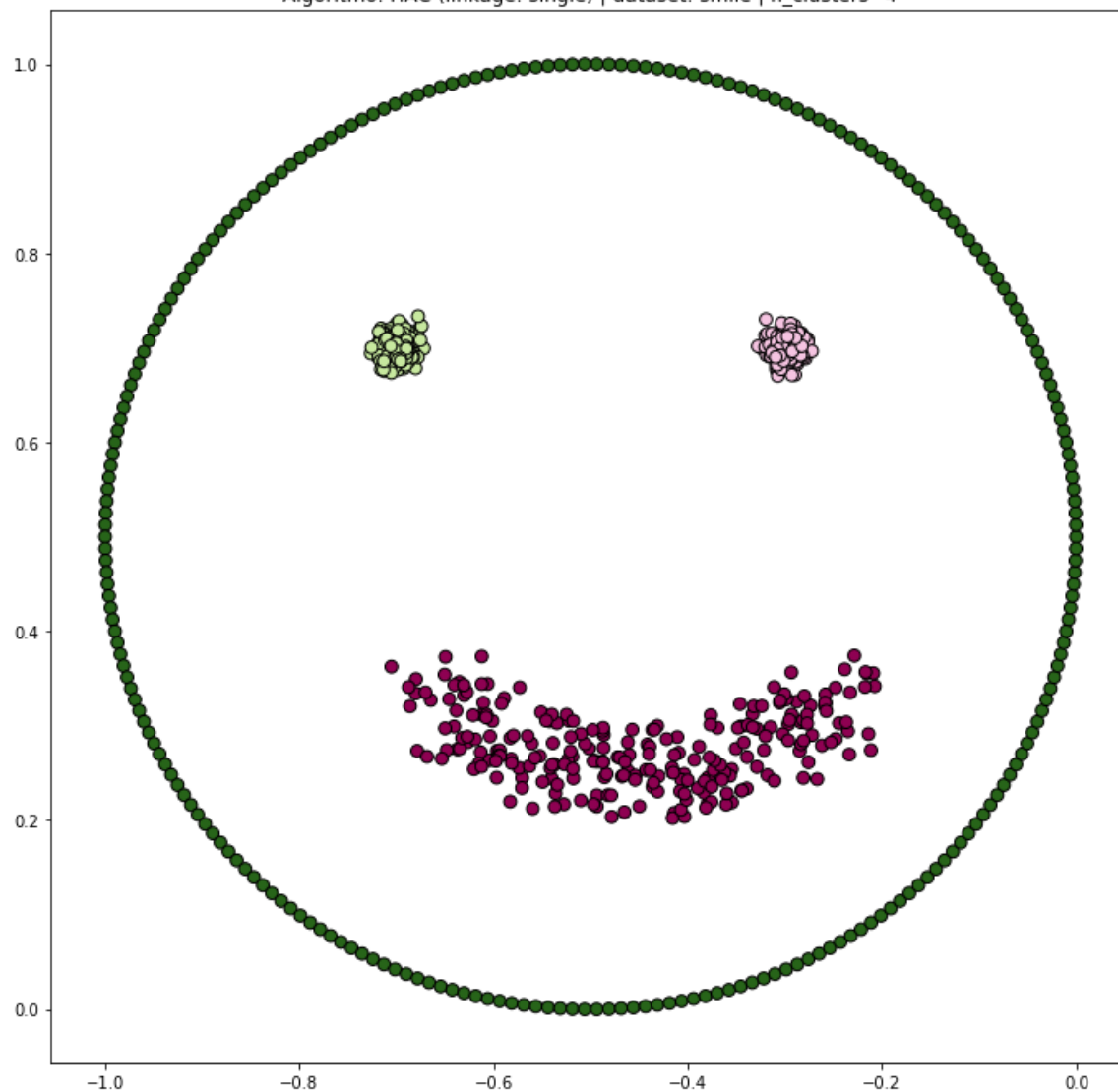
```
1 from sklearn.cluster import AgglomerativeClustering
2
3 def apply_hac(dataset, linkage, n_clusters):
4     hac = AgglomerativeClustering(n_clusters=n_clusters, linkage=linkage)
5     hac.fit(dataset)
6     return hac.labels_
7
8 @interact(
9     dataset_name=['smile', 'mouse', 'spiral'],
10    linkage=['single', 'complete'],
11    n_clusters=(2,10, 1),
12 )
13 def plot_hac(dataset_name, linkage, n_clusters):
14     if dataset_name == 'smile':
15         dataset = smile
16     elif dataset_name == 'mouse':
17         dataset = mouse
18     elif dataset_name == 'spiral':
19         dataset = spiral
20
21     labels = apply_hac(dataset, linkage, n_clusters)
22     plt.figure(figsize=(12,12))
23     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
24                 edgecolors='k', s=60, cmap=plt.cm.PiYG)
25     plt.title('Algoritmo: HAC (linkage: {}) | dataset: {} | n_clusters={}'.format(linkage, dataset_name, n_clusters))
26
```

dataset_na...

linkage

n_clusters

Algoritmo: HAC (linkage: single) | dataset: smile | n_clusters=4



Análisis Agglomerative Hierarchical Clustering

Bla...

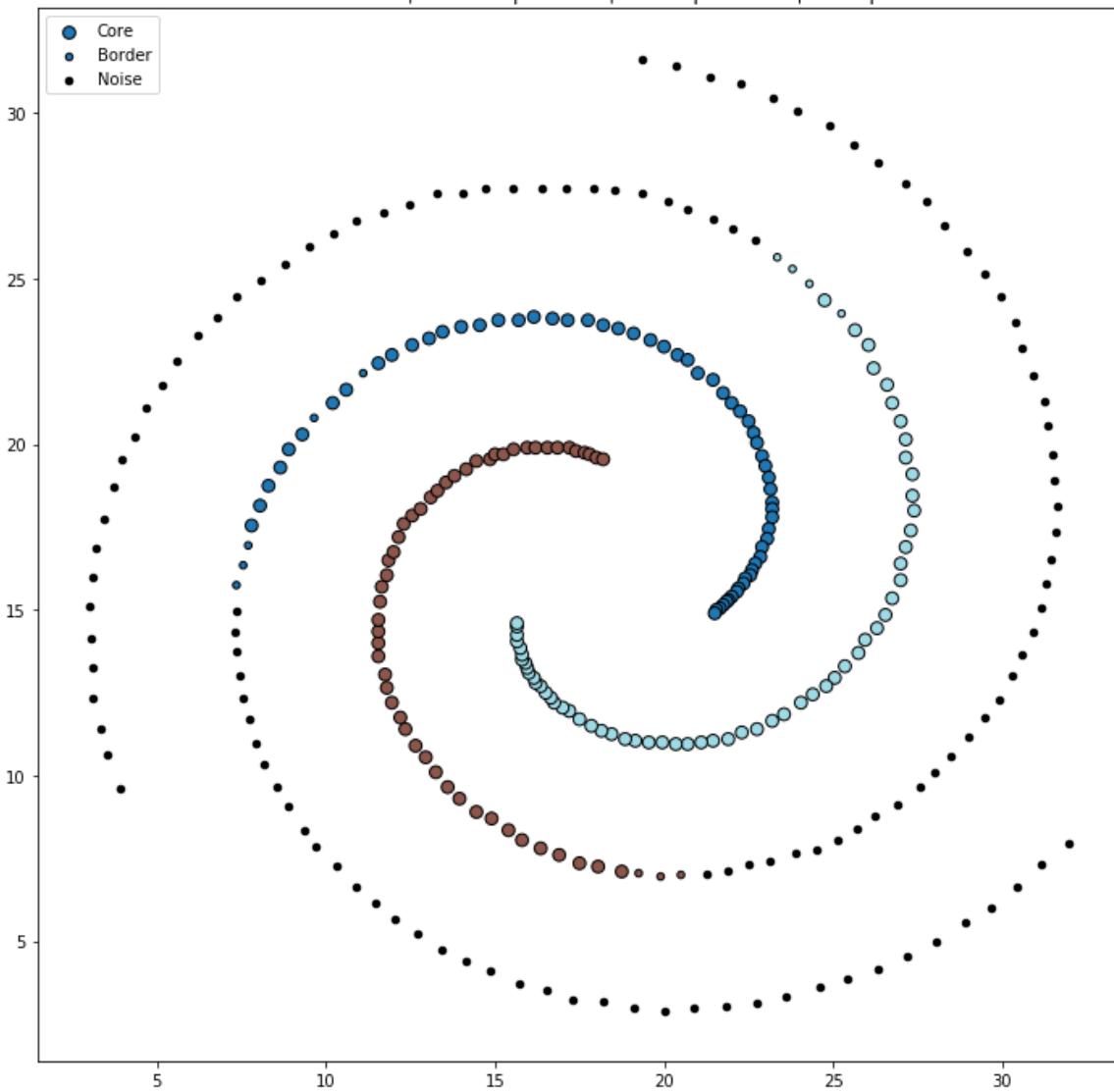
3. DBSCAN

In [62]:

```
1 from sklearn.cluster import DBSCAN
2
3 def apply_dbscan(dataset, min_pts, eps):
4     dbscan = DBSCAN(eps=eps, min_samples=min_pts)
5     dbscan.fit(dataset)
6     core_samples_mask = np.zeros_like(dbscan.labels_, dtype=bool)
7     core_samples_mask[dbscan.core_sample_indices_] = True
8     noise_points_mask = (dbscan.labels_ == -1)
9     border_points_mask = np.zeros_like(dbscan.labels_, dtype=bool)
10    border_points_mask[~core_samples_mask & ~noise_points_mask] = True
11
12    # Number of clusters in labels, ignoring noise if present.
13    n_clusters_ = len(set(dbscan.labels_)) - (1 if -1 in dbscan.labels_ \
14                                              else 0)
15
16    n_noise_ = list(dbscan.labels_).count(-1)
17    return dbscan.labels_, n_clusters_, n_noise_, core_samples_mask, \
18           border_points_mask, noise_points_mask
19
20 @interact(
21     dataset_name=['smile', 'mouse', 'spiral'],
22     min_pts=(1,50, 1),
23     eps=(0.01, 5.0, 0.01),
24 )
25 def plot_dbscan(dataset_name, min_pts, eps):
26     if dataset_name == 'smile':
27         dataset = smile
28     elif dataset_name == 'mouse':
29         dataset = mouse
30     elif dataset_name == 'spiral':
31         dataset = spiral
32
33     labels, n_clusters, n_noise, core_samples_mask, border_points_mask, \
34     noise_points_mask = apply_dbscan(dataset, min_pts, eps)
35     core_points = dataset[core_samples_mask]
36     border_points = dataset[border_points_mask]
37     noise_points = dataset[noise_points_mask]
38
39     n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
40     plt.figure(figsize=(12,12))
41
42     # Plot core samples
43     plt.scatter(core_points[:, 0], core_points[:, 1], marker='o',
44                 c=labels[core_samples_mask], edgecolors='k', s=60,
45                 cmap=plt.cm.tab20, label='Core')
46
47     # Plot border points
48     plt.scatter(border_points[:, 0], border_points[:, 1], marker='o',
49                 c=labels[border_points_mask], edgecolors='k', s=20,
50                 cmap=plt.cm.tab20, label='Border')
51
52     # Plot noise points
53     plt.scatter(noise_points[:, 0], noise_points[:, 1], marker='o',
54                 c='black', edgecolors='k', s=20,
55                 cmap=plt.cm.tab20, label='Noise')
56
57     plt.title('Algoritmo: DBSCAN | dataset: {} | eps={} | min_samples={} \n\
58 Clusters resultantes: {} | Core samples: {} | Border points: {} | Noise points: {}'.
59               .format(dataset_name, eps, min_pts, n_clusters,
60                       len(core_points), len(border_points),
61                       len(noise_points)))
62     plt.legend(loc='upper left')
```

min_pts 7
eps 1.97

Algoritmo: DBSCAN | dataset: spiral | eps=1.97 | min_samples=7
Clusters resultantes: 3 | Core samples: 174 | Border points: 12 | Noise points: 126



Análisis DBSCAN

Bla...

4. Mean-shift

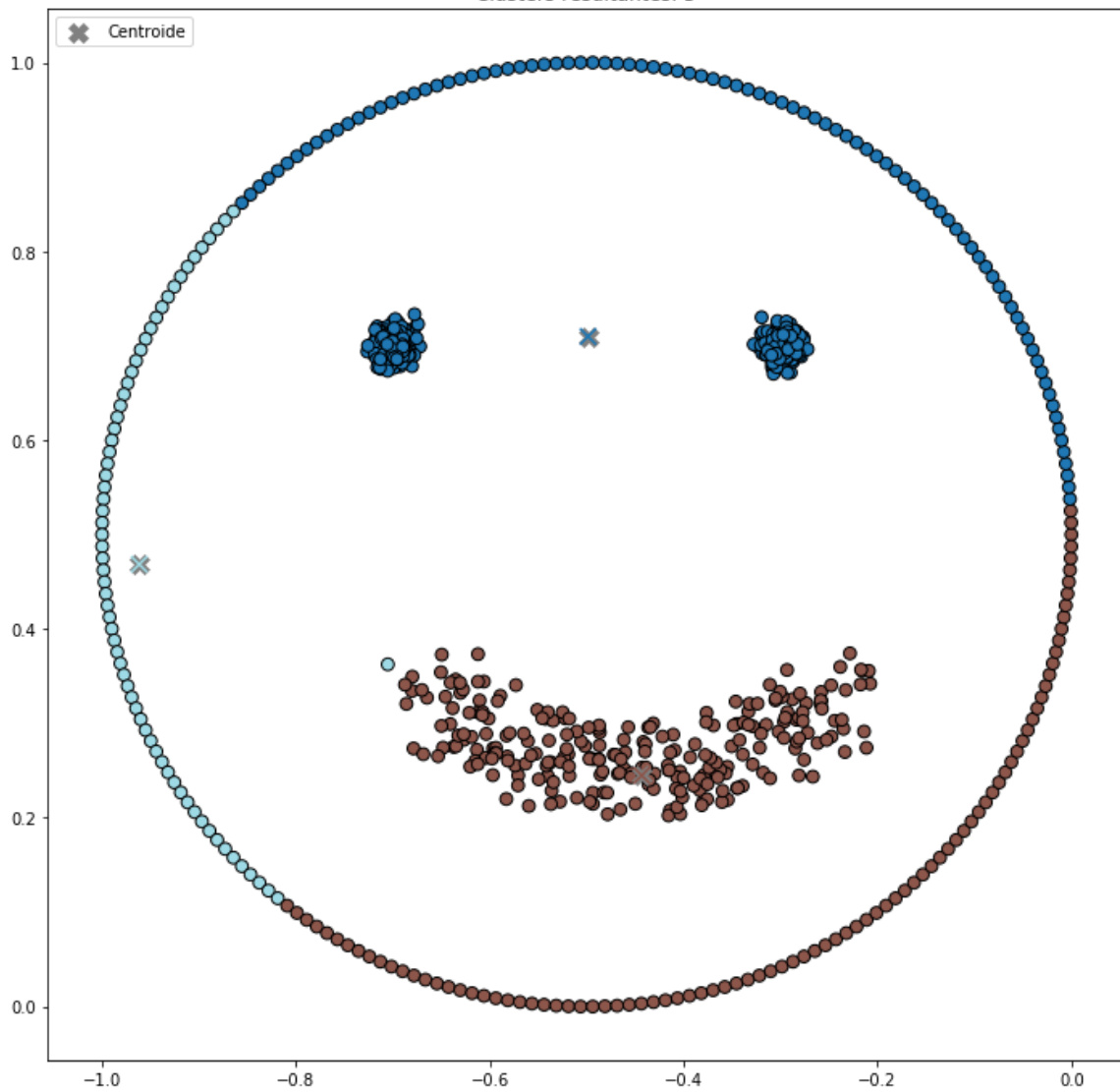
In [67]:

```
1 from sklearn.cluster import MeanShift
2
3 def apply_meanshift(dataset, bandwidth):
4     meanshift = MeanShift(bandwidth=bandwidth)
5     meanshift.fit(dataset)
6     return meanshift.cluster_centers_, meanshift.labels_
7
8 @interact(
9     dataset_name=['smile', 'mouse', 'spiral'],
10    bandwidth=(0.1, 10, 0.1),
11 )
12 def plot_kmeans(dataset_name, bandwidth):
13     if dataset_name == 'smile':
14         dataset = smile
15     elif dataset_name == 'mouse':
16         dataset = mouse
17     elif dataset_name == 'spiral':
18         dataset = spiral
19     centroids, labels = apply_meanshift(dataset, bandwidth=bandwidth)
20     plt.figure(figsize=(12,12))
21     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
22                 edgecolors='k', s=60, cmap=plt.cm.tab20)
23     plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=150,
24                 linewidths=.5, c='gray', cmap=plt.cm.tab20, label='Centroide')
25     plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=100,
26                 linewidths=2, c=list(range(len(centroids))),
27                 cmap=plt.cm.tab20)
28     plt.title('Algoritmo: Mean-shift | dataset: {} | bandwidth={} \n Clusters resultantes')
29     plt.legend(loc='upper left')
```

dataset_na...

bandwidth  0.30

Algoritmo: Mean-shift | dataset: smile | bandwidth=0.3
Clusters resultantes: 3



Análisis Mean-shift

Bla...

5. Spectral clustering

In [127]:

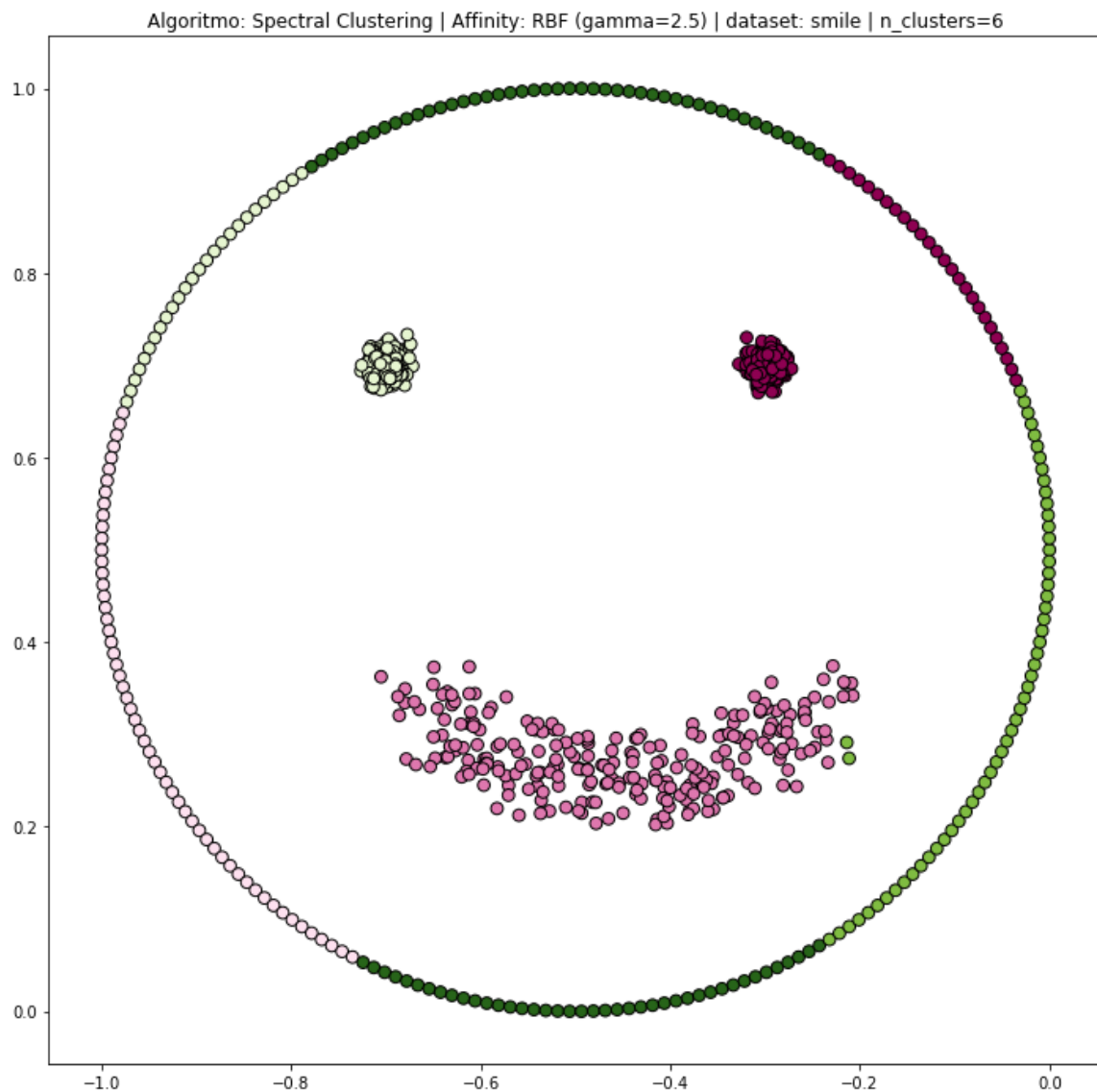
```
1  from sklearn.cluster import SpectralClustering
2
3  def apply_spectral_rbf(
4      dataset,
5      n_clusters,
6      random_state=0,
7      n_init=1,
8      gamma=1.0,
9      affinity_matrix_method='rbf',
10 ):
11     spectral = SpectralClustering(
12         n_clusters=n_clusters,
13         random_state=random_state,
14         n_init=n_init,
15         gamma=gamma,
16         affinity=affinity_matrix_method,
17     )
18     spectral.fit(dataset)
19     return spectral.labels_
20
21 def apply_spectral_nearest_neighbors(
22     dataset,
23     n_clusters,
24     n_neighbors,
25     random_state=0,
26     n_init=1,
27     affinity_matrix_method='nearest_neighbors',
28 ):
29     spectral = SpectralClustering(
30         n_clusters=n_clusters,
31         random_state=random_state,
32         n_init=n_init,
33         affinity=affinity_matrix_method,
34         n_neighbors=n_neighbors,
35     )
36     spectral.fit(dataset)
37     return spectral.labels_
38
39 @interact(
40     dataset_name=['smile', 'mouse', 'spiral'],
41     affinity=['RBF', 'K-nearest neighbors'],
42     n_clusters=(2,10, 1),
43     gamma=(0.1, 5.0, 0.1),
44     n_neighbors=(1, 20, 1),
45 )
46 def plot_spectral(
47     dataset_name,
48     affinity,
49     n_clusters,
50     gamma,
51     n_neighbors,
52 ):
53     if dataset_name == 'smile':
54         dataset = smile
55     elif dataset_name == 'mouse':
56         dataset = mouse
57     elif dataset_name == 'spiral':
58         dataset = spiral
59     if affinity == 'RBF':
```

```

60     labels = apply_spectral_rbf(
61         dataset=dataset,
62         n_clusters=n_clusters,
63         gamma=gamma,
64         affinity_matrix_method='rbf',
65     )
66     elif affinity == 'K-nearest neighbors':
67         labels = apply_spectral_nearest_neighbors(
68             dataset=dataset,
69             n_clusters=n_clusters,
70             n_neighbors=n_neighbors,
71             affinity_matrix_method='nearest_neighbors',
72         )
73     plt.figure(figsize=(12,12))
74     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
75                 edgecolors='k', s=60, cmap=plt.cm.PiYG)
76     if affinity == 'RBF':
77         plt.title('Algoritmo: Spectral Clustering | Affinity: RBF\
78 (gamma={}) | dataset: {} | n_clusters={}'.format(round(gamma, 2),
79                                                     dataset_name, n_clusters))
80     elif affinity == 'K-nearest neighbors':
81         plt.title('Algoritmo: Spectral Clustering | Affinity: {}-nearest\
82 neighbors | dataset: {} | n_clusters={}'.format(n_neighbors, dataset_name,
83                                                     n_clusters))
84

```

dataset_na...	<input type="text" value="smile"/>
affinity	<input type="text" value="RBF"/>
n_clusters	<input type="range" value="6"/>
gamma	<input type="range" value="2.50"/>
n_neighbors	<input type="range" value="10"/>



(Se pueden agregar más métodos de construcción de matriz de afinidad)

Análisis Spectral clustering

Bla...

6. Fuzzy C-Means

In [145]:

```
1 from skfuzzy import cluster as fuzzy
2 from ipywidgets import IntSlider
3
4 # Info: https://pythonhosted.org/scikit-fuzzy/auto\_examples/plot\_cmeans.html
5 def apply_cmeans(dataset, c, m, error, maxiter, seed=0):
6     # El argumento 'data' de cmeans exige que el dataset venga transpuesto!
7     cntr, u, u0, d, jm, p, fpc = fuzzy.cmeans(
8         data=dataset.T,
9         c=c,
10        m=m,
11        error=error,
12        maxiter=maxiter,
13        seed=seed,
14    )
15    return cntr, u, u0, d, jm, p, fpc
16
17 @interact(
18     dataset_name=['smile', 'mouse', 'spiral'],
19     n_clusters=(2, 10, 1),
20     p_exponent=(1.1, 3.0, 0.1),
21     error=FloatSlider(
22         min=5e-4,
23         max=5e-2,
24         step=5e-4,
25         continuous_update=False,
26         readout=True,
27         readout_format='.4f'
28     ),
29     max_iterations=(1, 1000, 1),
30 )
31 def plot_cmeans(
32     dataset_name,
33     n_clusters,
34     p_exponent,
35     error,
36     max_iterations,
37 ):
38     if dataset_name == 'smile':
39         dataset = smile
40     elif dataset_name == 'mouse':
41         dataset = mouse
42     elif dataset_name == 'spiral':
43         dataset = spiral
44     cntr, u, u0, d, jm, p, fpc = apply_cmeans(
45         dataset=dataset,
46         c=n_clusters,
47         m=p_exponent,
48         error=error,
49         maxiter=max_iterations,
50     )
51     # El color de cada punto es asignado segun el cluster al cual pertenezca
52     # con mayor porcentaje
53     labels = u.argmax(axis=0)
54     plt.figure(figsize=(12,12))
55     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
56                 edgecolors='k', s=60, cmap=plt.cm.ocean)
57     plt.scatter(cntr[:, 0], cntr[:, 1], marker='x', s=150,
58                 linewidths=.5, c='gray', cmap=plt.cm.ocean, label='Centroide')
59     plt.scatter(cntr[:, 0], cntr[:, 1], marker='x', s=100,
```



```

60         linewidths=2, c=list(range(n_clusters)),
61         cmap=plt.cm.ocean)
62     plt.title('Algoritmo: Fuzzy C-Means | dataset: {} | n_clusters={} | \
63 p={} | Máx. Iters={} | error={} \n Fuzzy partition coefficient (FPC): {}'.format(
64         dataset_name,
65         n_clusters,
66         round(p_exponent, 2),
67         max_iterations,
68         round(error, 5),
69         round(fpc, 3),
70     ))
71     plt.legend(loc='upper left')

```

dataset_na...

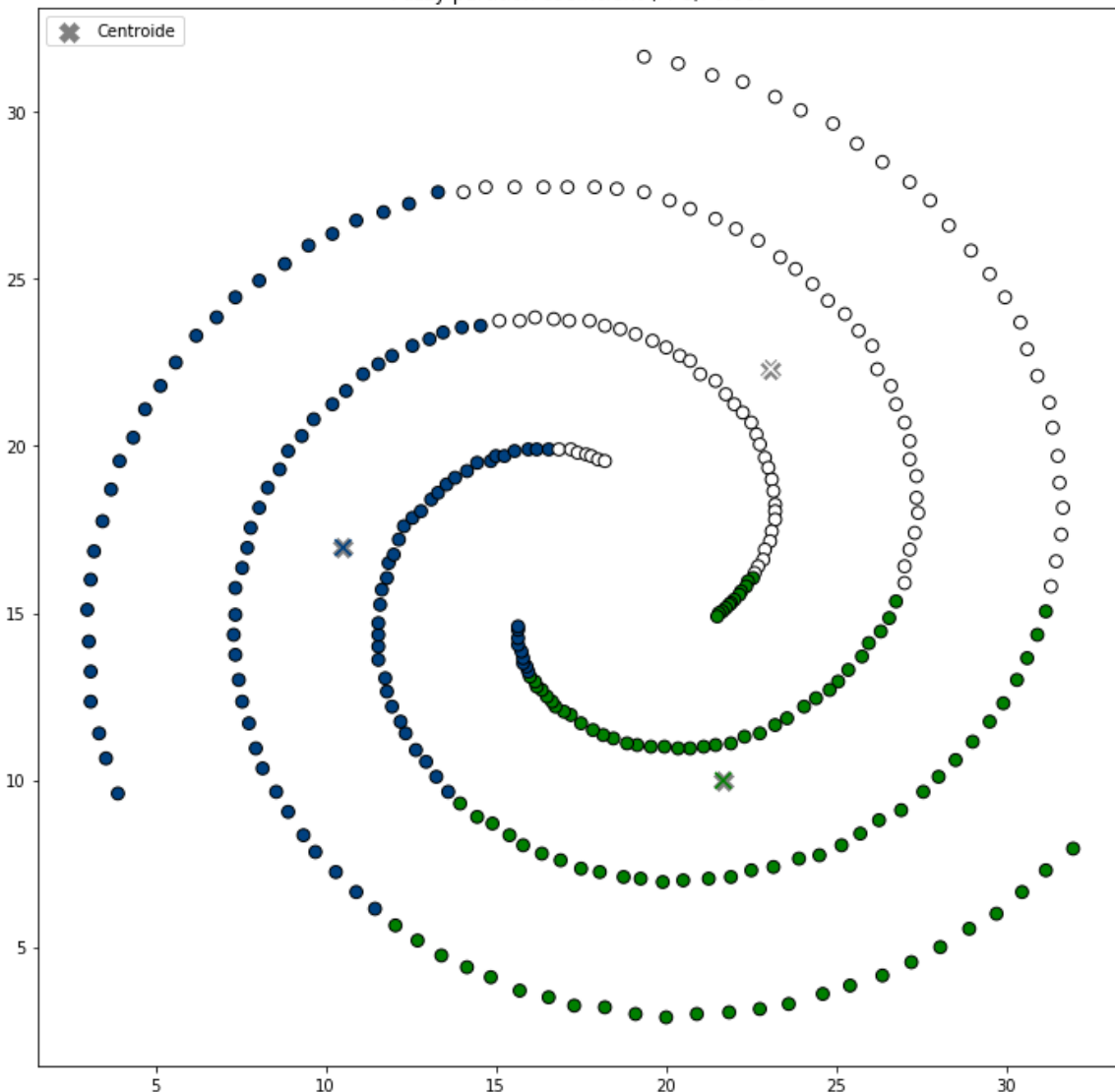
n_clusters 3

p_exponent 2.00

error 0.0015

max_iterati... 1000

Algoritmo: Fuzzy C-Means | dataset: spiral | n_clusters=3 | p=2.0 | Máx. Iters=1000 | error=0.0015
Fuzzy partition coefficient (FPC): 0.603



(este link puede ayudar con el análisis: https://pythonhosted.org/scikit-fuzzy/auto_examples/plot_cmeans.html (https://pythonhosted.org/scikit-fuzzy/auto_examples/plot_cmeans.html))

Bla...

Parte II

(a) Se tiene un conjunto de datos con 100 objetos. Se le pide realizar clustering utilizando K-means, pero para todos los valores de k , $1 \leq k \leq 100$, el algoritmo retorna que todos los clusters estan vacíos, excepto uno. ¿En que situación podría ocurrir esto? (analice los datos y no los parametros del algoritmo, i.e., iteraciones). ¿Qué resultado tendría single-link y DBSCAN para este tipo de datos?

1 Resp. :

(b) Considerando single-link y complete-link hierarchical clustering, ¿es posible que un objeto esté más cerca (en distancia Euclidiana) de los objetos de otros clusters en relación a los de su propio cluster? Si fuese posible, ¿en que enfoque (single y/o complete) esto podría ocurrir? Justifique con un ejemplo en cada caso.

Resp.: En cada iteración de single-link, un punto se incluye en cualquiera sea el cluster al cual pertenezca su punto más cercano. Al construirse los clusters de esta manera, no existe forma de que algún par de puntos de distintos clusters sean más cercanos entre sí que con los puntos de sus respectivos clusters.

En cambio, en complete-link esto sí puede ocurrir. Considerar, por ejemplo, la siguiente distribución de los objetos A, B, C, D, E en una dimensión:

A-B-C--D---E

Donde la distancia es directamente proporcional a la cantidad de caracteres entre los objetos. Si queremos un resultado de 2 clusters, se obtiene:

[A-B-C]--[D---E]

De aquí notamos que D dista de C (que está en un cluster distinto) sólo en 2 unidades, mientras que la distancia con E (que está en el mismo cluster) es de 3. Así, con este ejemplo, queda demostrada la posibilidad de que un objeto de un cluster puede tener como objetos más cercanos a aquellos de otros clusters.

In []:

1