



# 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')
```

**(Hacer plot y breve análisis de cada dataset: hablar sobre cantidad de datos, presencia obvia de clusters, densidad de éstos, convexidad, etc.)**

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

### 1. K-Means

```

In [64]: 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(min=5e-5, max=5e-4, step=5e-5, continuous_update=False),
23 )
24 def plot_kmeans(dataset_name, k, max_iterations, tolerance):
25     if dataset_name == 'smile':
26         dataset = smile
27     elif dataset_name == 'mouse':
28         dataset = mouse
29     elif dataset_name == 'spiral':
30         dataset = spiral
31     centroids, labels = apply_kmeans(dataset, k, max_iterations, tolerance)
32     plt.figure(figsize=(12,12))
33     plt.scatter(dataset[:, 0], dataset[:, 1], marker='o', c=labels,
34                edgecolors='k', s=60, cmap=plt.cm.ocean)
35     plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=150,
36                linewidths=.5, c='gray', cmap=plt.cm.ocean, label='Centroide')
37     plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=100,
38                linewidths=2, c=list(range(len(centroids))),
39                cmap=plt.cm.ocean)
40     plt.title('Algoritmo: KMeans | dataset: {} | k={} | Máx. Iters={} | Tolerancia={}'.format(dataset_name, k, max_iterations, tolerance))
41     plt.legend(loc='upper left')

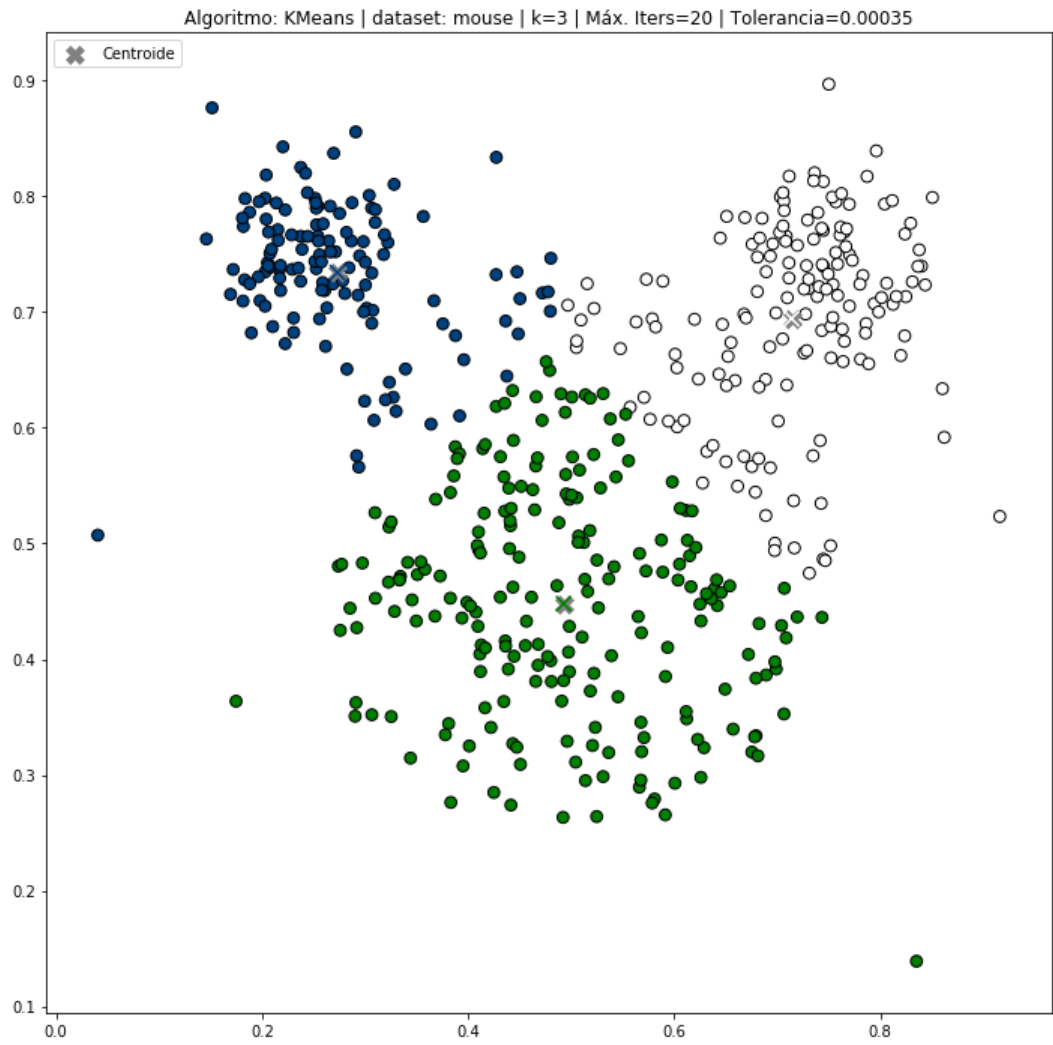
```

dataset\_na...

k  3

max\_iterati...  20

tolerance  0.00



**Análisis K-Means**

Bla...

**2. Agglomerative Hierarchical Clustering**

```

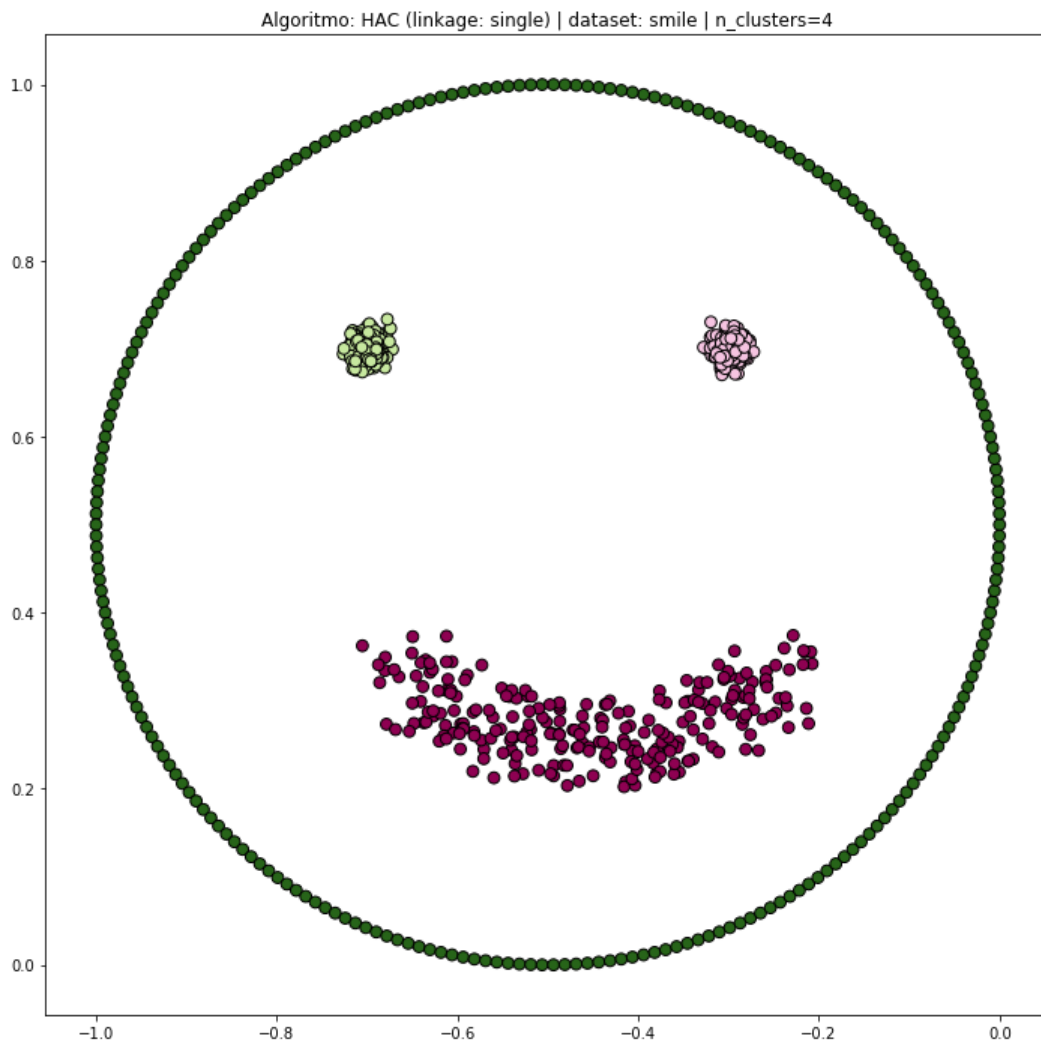
In [42]: 1 from sklearn.cluster import AgglomerativeClustering
2
3 def apply_hac(dataset, linkage, n_clusters):
4     hac = AgglomerativeClustering(n_clusters=n_clusters, linkage=linkage)
5     hac.fit(dataset)
6     return hac.labels_
7
8 @interact(
9     dataset_name=['smile', 'mouse', 'spiral'],
10    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,
26

```

dataset\_na...

linkage

n\_clusters



### Análisis Agglomerative Hierarchical Clustering

Bla...

**3. DBSCAN**

```

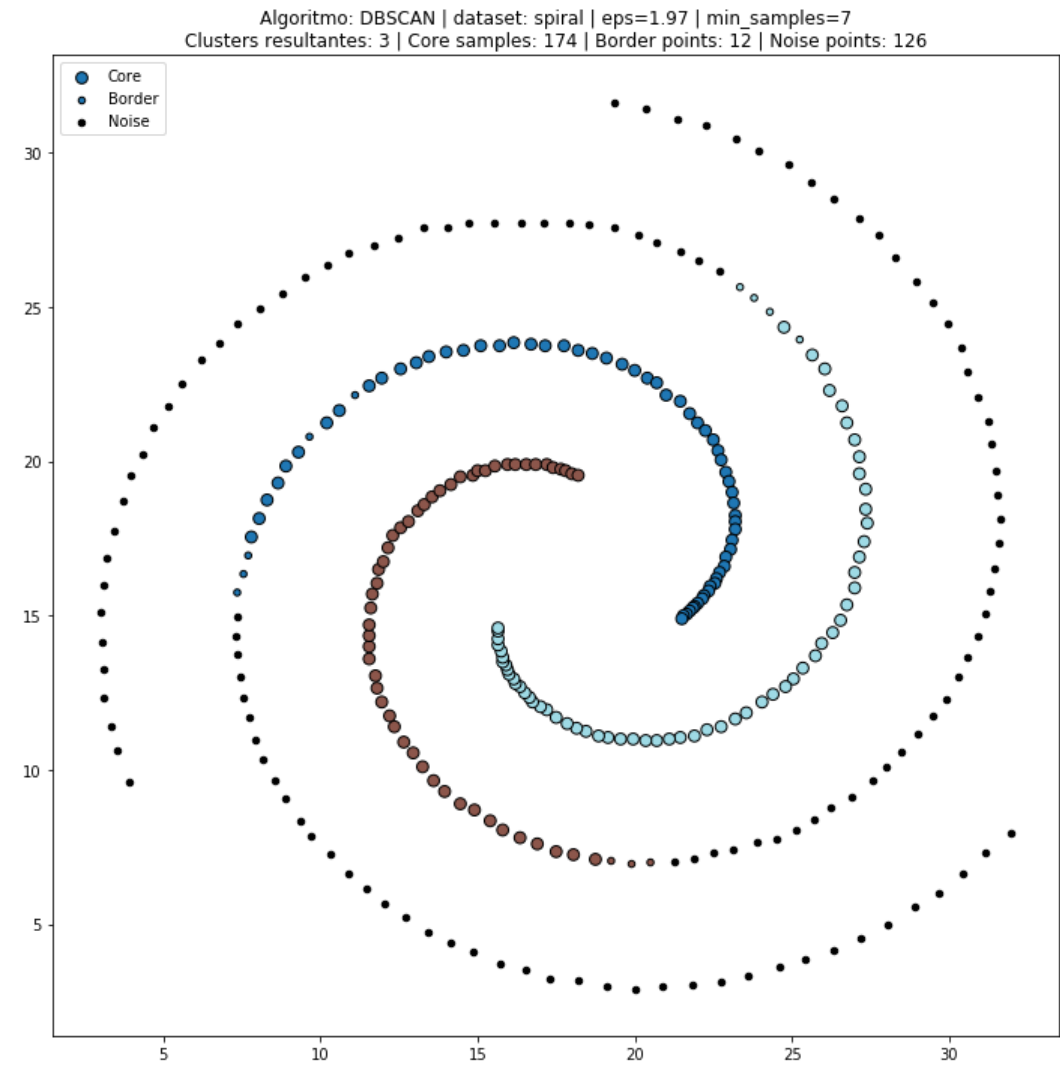
In [62]: 1 from sklearn.cluster import DBSCAN
2
3 def apply_dbscan(dataset, min_pts, eps):
4     dbscan = DBSCAN(eps=eps, min_samples=min_pts)
5     dbscan.fit(dataset)
6     core_samples_mask = np.zeros_like(dbscan.labels_, dtype=bool)
7     core_samples_mask[dbscan.core_sample_indices_] = True
8     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     n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
39     plt.figure(figsize=(12,12))
40     # Plot core samples
41     plt.scatter(core_points[:, 0], core_points[:, 1], marker='o',
42                c=labels[core_samples_mask], edgecolors='k', s=60,
43                cmap=plt.cm.tab20, label='Core')
44     # Plot border points
45     plt.scatter(border_points[:, 0], border_points[:, 1], marker='o',
46                c=labels[border_points_mask], edgecolors='k', s=20,
47                cmap=plt.cm.tab20, label='Border')
48     # Plot noise points
49     plt.scatter(noise_points[:, 0], noise_points[:, 1], marker='o',
50                c='black', edgecolors='k', s=20,
51                cmap=plt.cm.tab20, label='Noise')
52     plt.title('Algoritmo: DBSCAN | dataset: {} | eps={} | min_samples={}\n\
53     Clusters resultantes: {} | Core samples: {} | Border points: {} | Noise points: {}'.
54               .format(dataset_name, eps, min_pts, n_clusters,
55                       len(core_points), len(border_points),
56                       len(noise_points)))
57     plt.legend(loc='upper left')

```

dataset\_na... spiral

min\_pts

eps



**Análisis DBSCAN**

Bla...

**4. Mean-shift**

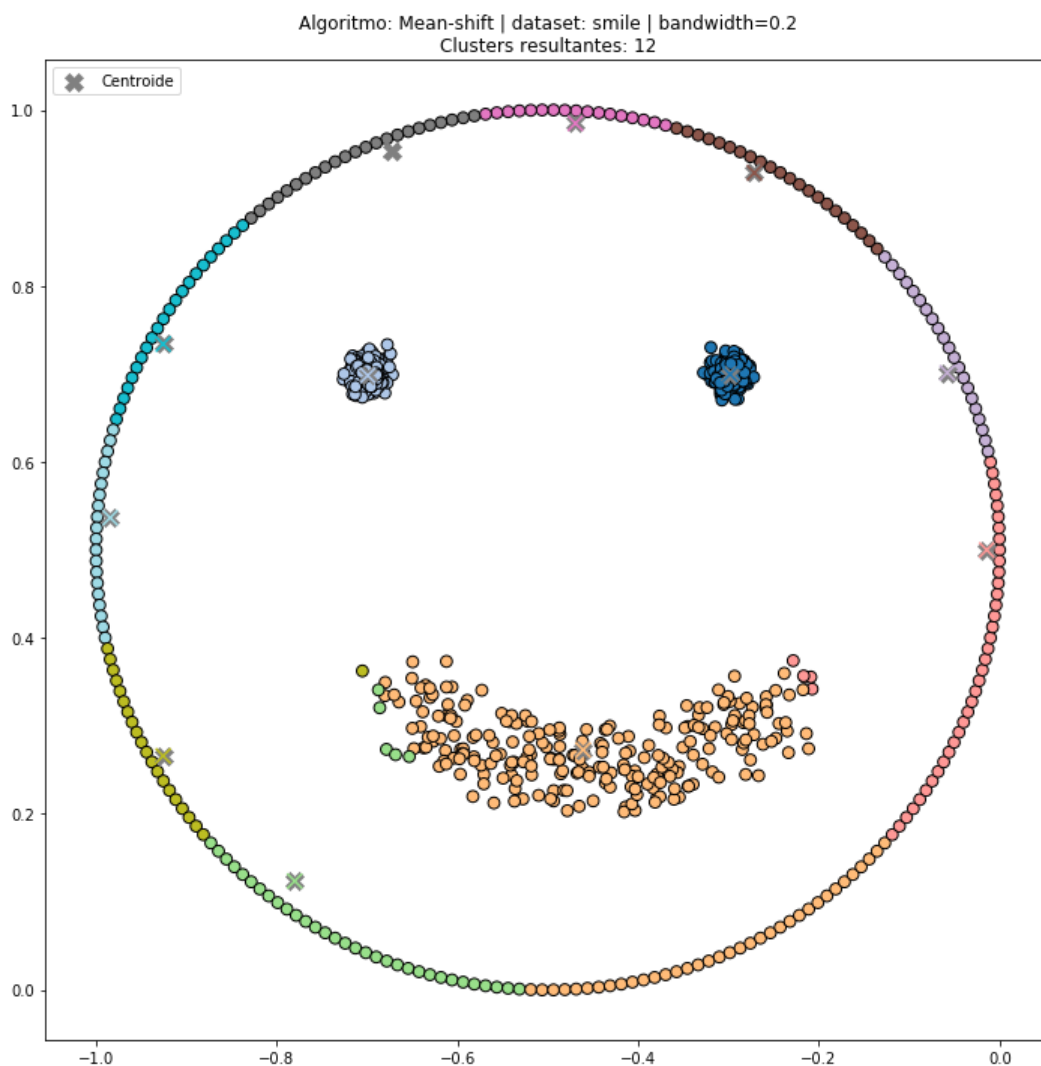
```

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={}\nClusters resultantes: {}'.format(datas
29     plt.legend(loc='upper left')

```

dataset\_na... smile

bandwidth 0.20



### Análisis Mean-shift



Bla...

## 5. Spectral clustering

```
In [ ]: 1 from sklearn.cluster import SpectralClustering
        2
        3 def generate_affinity_matrix(dataset, method):
        4     # METHOD -> PARAMETERS
        5     # epsilon-ball -> epsilon
        6     # k-nearest -> k
        7     # fully connected -> (no parameters)
        8     # RBF kernel -> delta (kernel width)
        9
        10 def apply_spectral(dataset, n_clusters, random_state=0, n_init=1,):
        11     return
```