



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')
```

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

1. K-Means

```

In [5]: 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, 300, 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], c=labels)
34     plt.scatter(centroids[:, 0], centroids[:, 1], marker='X', s=150, linewidths=.5, c='gray')
35     plt.scatter(centroids[:, 0], centroids[:, 1], marker='x', s=100, linewidths=2, c=list(range(len(centroids))))
36     plt.title('Algoritmo: KMeans | dataset: {} | k={} | Máx. Iters={} | Tolerancia={}'.format(dataset_name, k, max_iterations, tolerance))
37

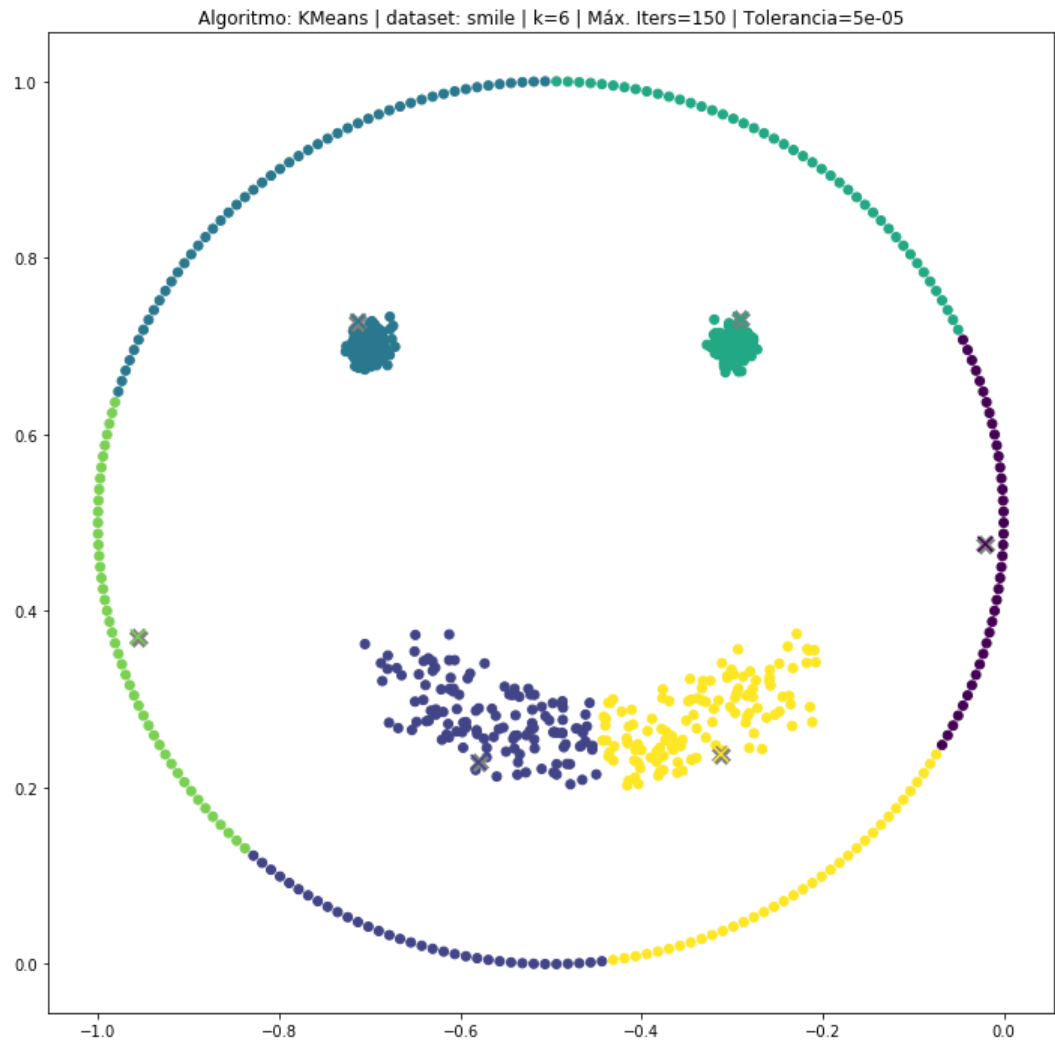
```

dataset_na... smile

k 6

max_iterati... 150

tolerance 0.00



Análisis K-Means

Bla...

2. Agglomerative Hierarchical Clustering

```

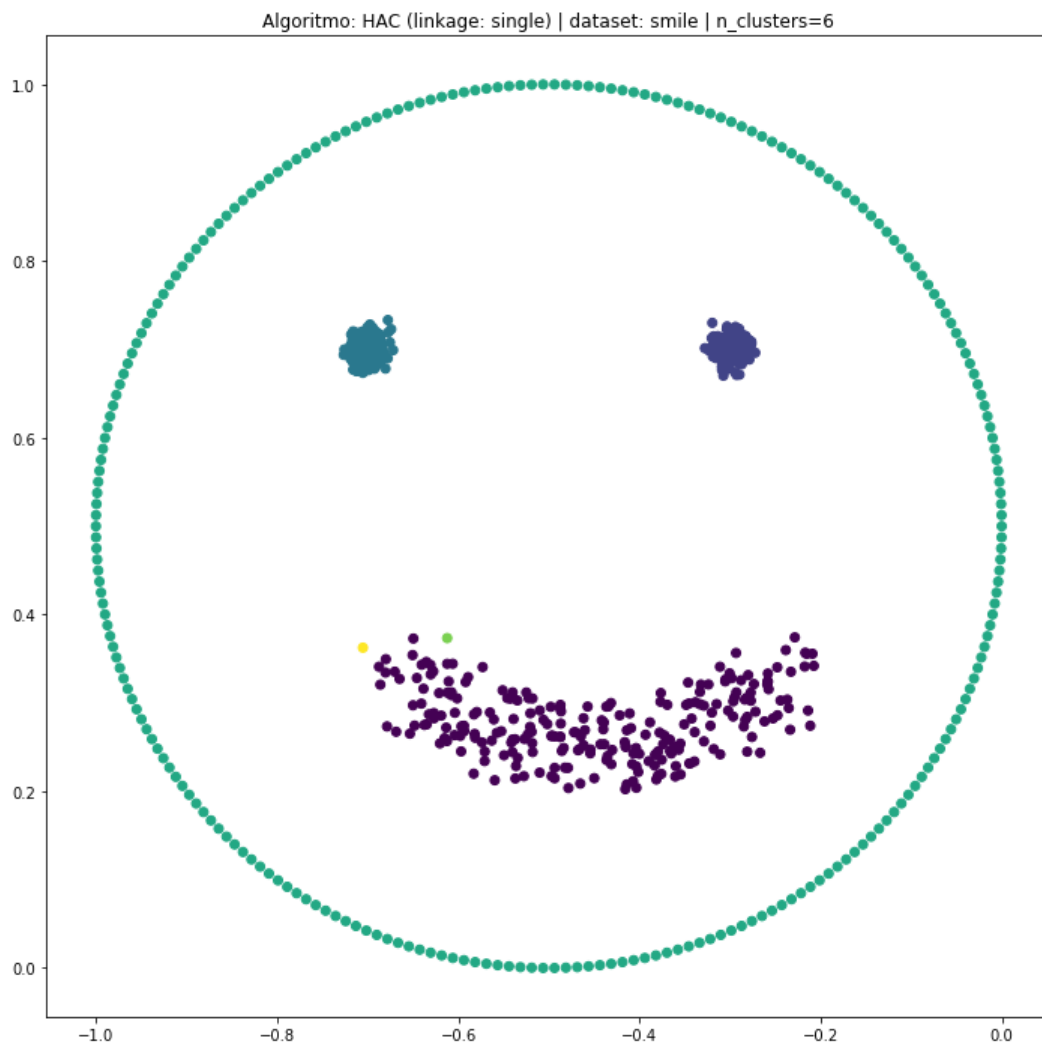
In [4]: 1 from sklearn.cluster import AgglomerativeClustering
2
3 def apply_hac(dataset, linkage, n_clusters):
4     hac = AgglomerativeClustering(n_clusters=n_clusters, linkage=linkage)
5     hac.fit(dataset)
6     return hac.labels_
7
8 @interact(
9     dataset_name=['smile', 'mouse', 'spiral'],
10    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], c=labels)
24     plt.title('Algoritmo: HAC (linkage: {}) | dataset: {} | n_clusters={}'.format(linkage, dataset_name,
25

```

dataset_na...

linkage

n_clusters



Análisis Agglomerative Hierarchical Clustering

Bla...

In []:

1