



CS2601
INTELIGENCIA ARTIFICIAL

Proyecto 3: Clustering de información genética en tejidos humanos

Mauricio Pinto
Francesco Uccelli
Juan Manuel Navarro

20 de Junio del 2021

1 Introducción

En este proyecto se utilizarán diferentes métodos de clustering para encontrar tejidos en una base de datos que contiene información genérica humana. Dado que la dimensionalidad de los datos es demasiado alta para trabajarla con los métodos directamente, se propone un método de reducción de dimensionalidad mediante análisis de componentes principales. Para el clustering se usarán los métodos de Gaussian Mixture Model, K-means, DBSCAN y Agglomerative Hierarchical Clustering. Finalmente se realizaran experimentos y se compararán los resultados.

2 Análisis

2.1 Datos

El dataset consiste en observaciones de medición de presencia de genes en tejidos humanos. En particular esta data representa expresiones de RNA 8 tejidos distintos, cada uno con múltiples individuos. Esta data es altamente dimensional lo cual dificulta mucho el análisis de su información bruta. Por esto se utiliza el método de análisis de componentes principales que se expone a continuación.

2.2 Componentes Principales

El método de reducción de dimensionalidad por análisis de componentes principales consiste en representar la información que brindan las variables en altas dimensiones en dimensiones bajas asumiendo una pérdida de información lo más baja posible. A continuación mostraremos como definimos este método. Para el algoritmo usamos PCA del módulo *sklearn.decomposition*.

A continuación podemos ver el gráfico que nos muestra con que cantidad de componentes se puede explicar que porcentaje de la varianza de los datos:

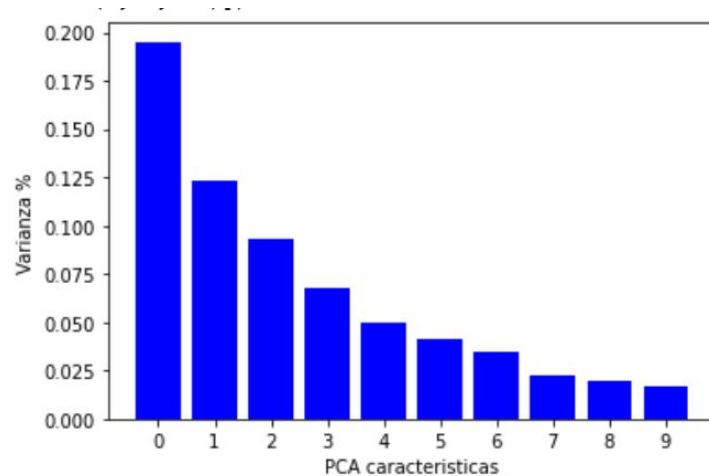


Figure 1: Features del PCA vs Varianza

Como podemos ver en la imagen, es suficiente un con 3 componentes (desde 0 en el gráfico) para representar un poco más de 50% de la varianza en los datos. Esto quiere decir que con esta cantidad de dimensiones deberíamos poder obtener resultados de clustering decentes. Además si mostramos la información que contienen los componentes 0 y 1 podemos ver claramente clusters definidos.

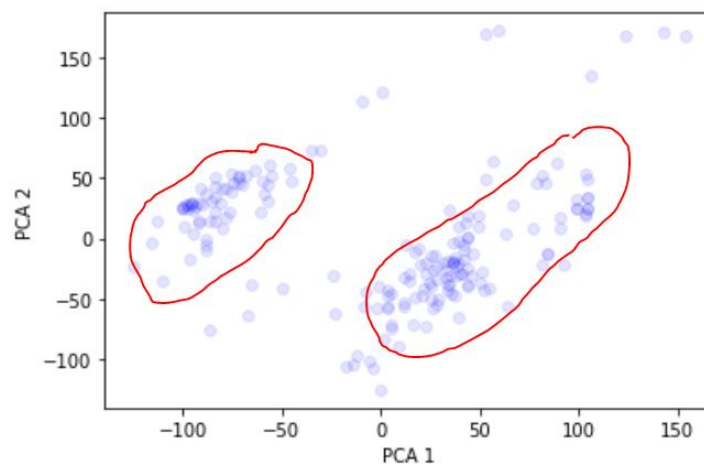


Figure 2: Features 1 y 2

Por otro lado, observemos también con que cantidad de clusters podemos explicar la data en su mayoría. En el siguiente gráfico se muestra la "inercia"

que tiene el significado de cuanto ayuda un nuevo cluster a la distribución de los datos y el número de clusters.

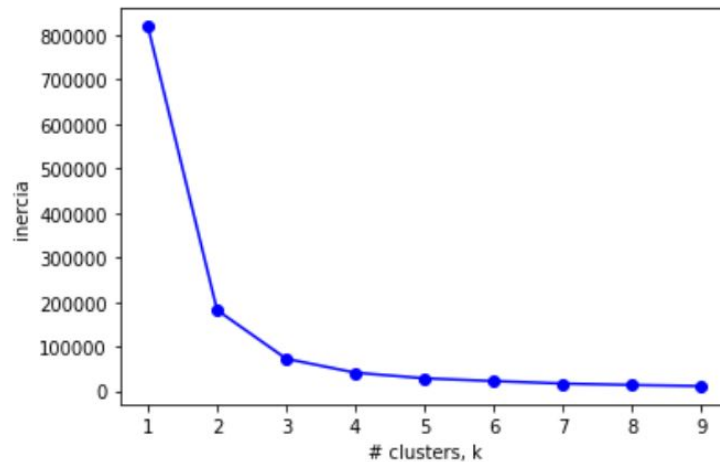


Figure 3: Inercia vs cantidad de clusters

Este gráfico explica que para una cantidad de clusters mayor a 4, aumentar un cluster más va a ayudar muy poco a explicar la distribución de los datos. Es por esto que, entonces, usando 3 dimensiones y 4 clusters se deberían de poder obtener resultados confiables.

3 Algoritmos Utilizados

3.1 K-Means

K-Means es un algoritmo de agrupación en clústeres no supervisado. Se usa cuando tenemos muchos datos sin etiquetar. El objetivo de este algoritmo es encontrar grupos "K" (clusters) entre los datos sin procesar.

El algoritmo funciona de forma iterativa para asignar a cada "punto" (las filas de nuestro conjunto de entrada forman una coordenada) uno de los grupos "K" en función de sus características, los cuales se agrupan según la similitud de tales características, devolviéndonos los centroides de cada grupo.

```
def classify (el, centers):  
    minn = 1000000  
    for key in centers:  
        if dist(el, centers[key]) < minn:  
            minn = dist(el, centers[key])  
            cluster = key  
    return cluster
```

```
def kmeans(df,k,cluster_centers):
    classes = {}
    for i in range (k):
        classes[str(i)] = []
    for x in range(len(df[0])):
        element = list ([df[0][x],df[1][x]])
        assigned_cluster =
            classify (element, cluster_centers)
        classes[assigned_cluster].append (element)
    return classes
```

3.2 GMM

Gaussian mixture model intenta encontrar una mezcla de distribuciones de probabilidad gaussianas multidimensionales que modelen mejor cualquier conjunto de datos de entrada. En el caso más simple, los GMM se pueden usar para encontrar clusters de la misma manera que k-means. Pero debido a que GMM contiene un modelo probabilístico bajo el capó, también es posible encontrar asignaciones de conglomerados probabilísticos.

```
def run_gmm (self, iterations, x, aux=None):
    means = self.gmm_generate_means ()
    prior = self.gmm_generate_prior ()
    covariance = self.gmm_generate_covariance (means, x)

    for i in range (iterations):
        gamma = self.gmm_calculate_gamma (prior,
            covariance, x, means)
        means = self.gmm_calculate_means (gamma, x)
        prior = self.gmm_calculate_prior (gamma, x)
        covariance = self.gmm_calculate_covariance (gamma,
            x, means)

    return means
```

3.3 DBSCAN

es un algoritmo básico para el agrupamiento basado en densidad. Puede descubrir grupos de diferentes formas y tamaños a partir de una gran cantidad de datos, que contienen ruido y valores atípicos.

K-means puede agrupar observaciones poco relacionadas. Cada observación eventualmente se convierte en parte de algún grupo, incluso si las observaciones están dispersas en el espacio vectorial. Un ligero cambio en los puntos de datos podría afectar el resultado de la agrupación. El DBSCAN reduce el problema debido a la forma en que se forman los clusters.

Una ventaja del DBSCAN es que no tienes que especificar el numero de clusters a utilizar, solo necesitas la función para calcular la distancia entre valores y definir la distancia que se pueda considerar cercana (minPoints y radius). DBSCAN también produce resultados mas razonables a comparación de K-means.

```
def dbscan(puntos, clusters, radius, tree):
    for i in range (len(puntos)):
        assigned_cluster = -1
        if clusters[i] == -1:
            assigned_cluster = i
        else:
            assigned_cluster = clusters[i]
        ind = tree.query_radius([puntos[i]], r=radius,
                                count_only=False, return_distance=False)
        is_cluster = False
        if len(ind[0]) >= 10:
            is_cluster = True
        if is_cluster:
            for m in ind[0]:
                clusters[m] = assigned_cluster
    return clusters
```

3.4 AHC

Agglomerative Hierarchical Clustering trae también ciertas ventajas sobre el resto. Funciona a partir de las diferencias entre los objetos a agrupar. Un tipo de disimilitud se puede adaptar al tema estudiado y a la naturaleza de los datos. Además, uno de los resultados es el dendrograma que muestra la agrupación progresiva de los datos. Entonces es posible tener una idea de un número adecuado de clases en las que se pueden agrupar los datos.

AHC funciona de manera iterativa, empezando por calcular la disimilitud entre N objetos. Luego juntando 2 objetos para ver si es que juntos minimizan un criterio de aglomeración, y finalmente calculando la disimilitud entre la clase y el resto de objetos utilizando el mismo criterio.

4 Experimentación

Para el caso de los gráficos se plantean diferentes cantidades de clusters en el k-means y en el DBSCAN los clusters se generan automáticamente. Por esto, puede resultar no completamente útil comparar medidas de precisión en casos donde la cantidad de clusters difiere. En cambio, planteamos a continuación el gráfico con los valores de las labels reales y comparamos los resultados de cada experimento visualmente.

Figure 4: Comparación entre diferentes algoritmos de clustering

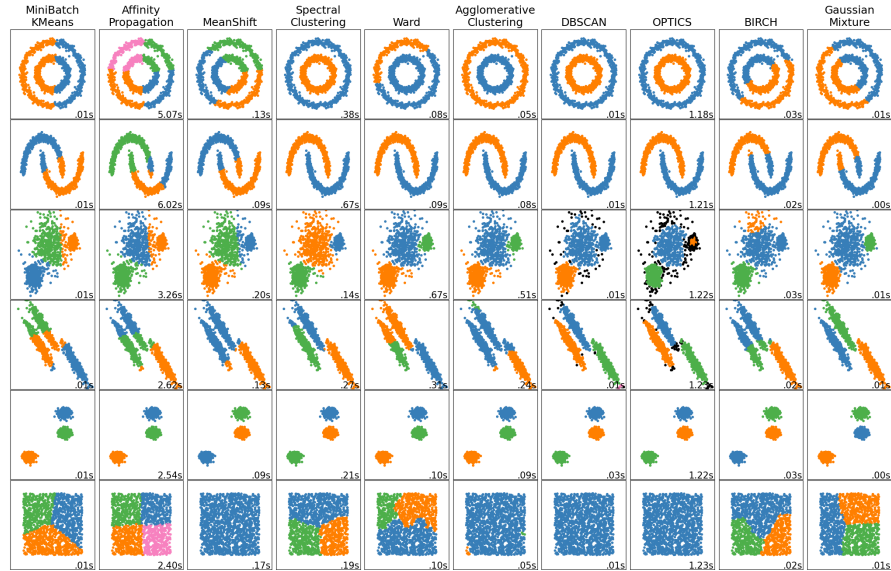
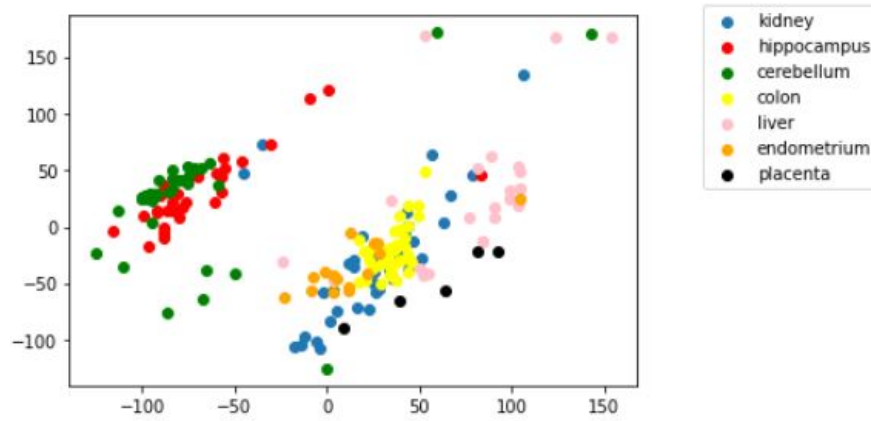


Figure 5: 3 clusters



4.1 K-Means

En este algoritmo variamos el parámetro k y mostramos los resultados

Figure 6: 3 clusters

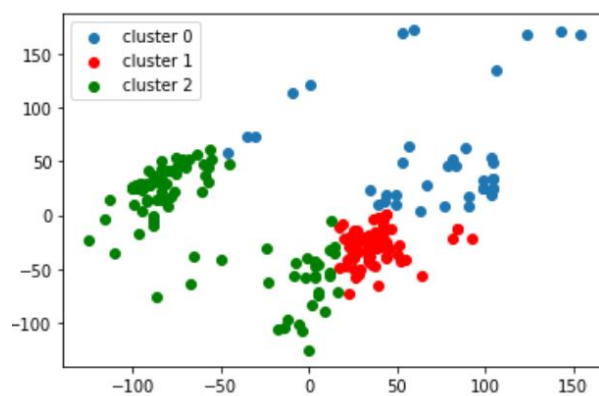


Figure 7: 5 clusters

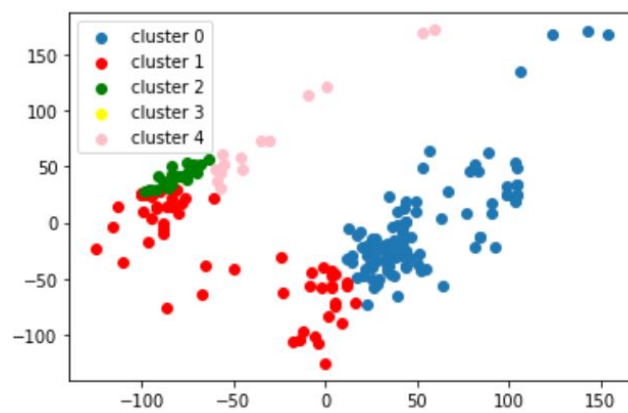
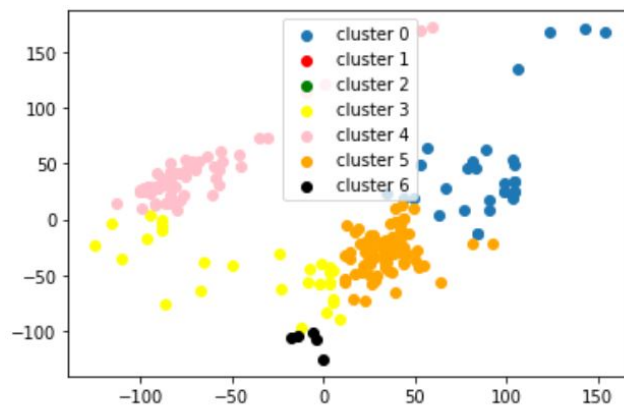


Figure 8: 7 clusters



4.2 DBSCAN

Figure 9: radio 10 minPoints 5

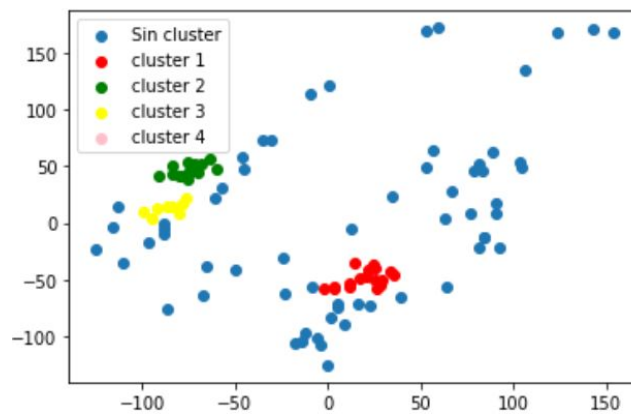


Figure 10: radio 20 minPoints 8

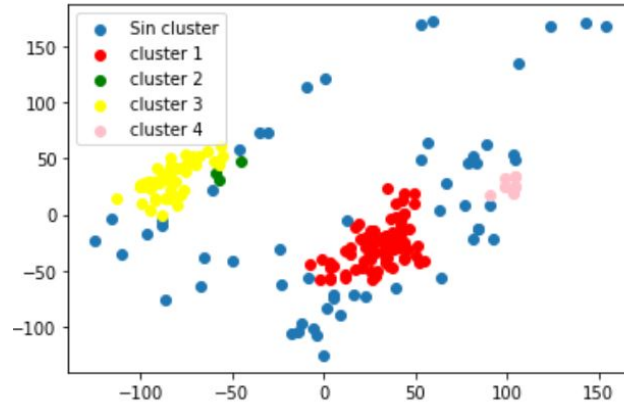
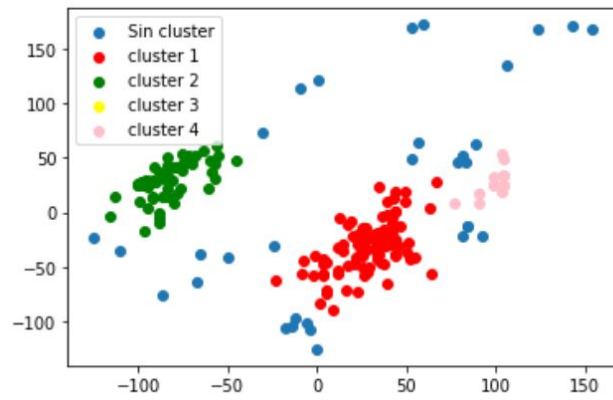


Figure 11: radio 10 minPoints 5



4.3 AHC

Figure 12: 5 clusters

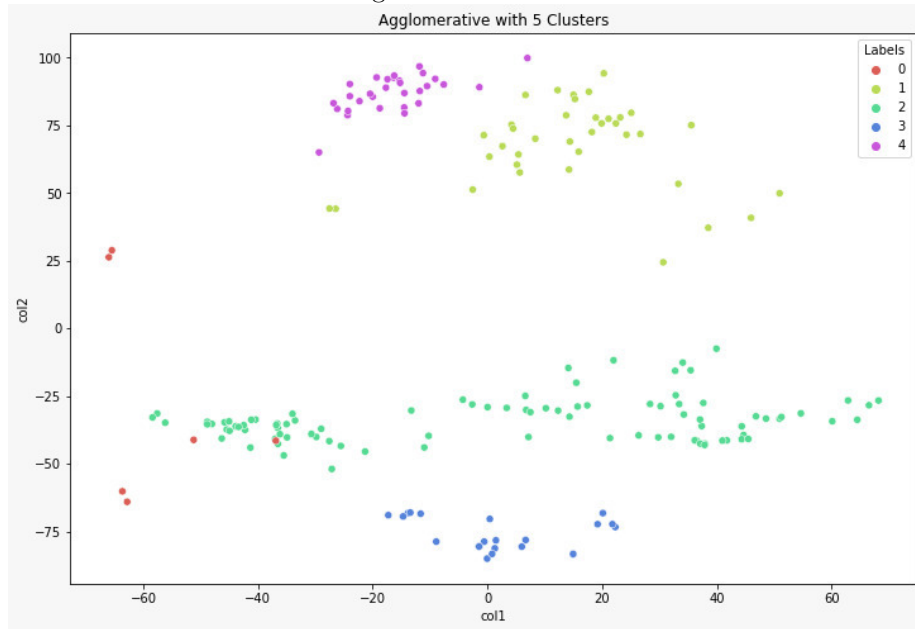
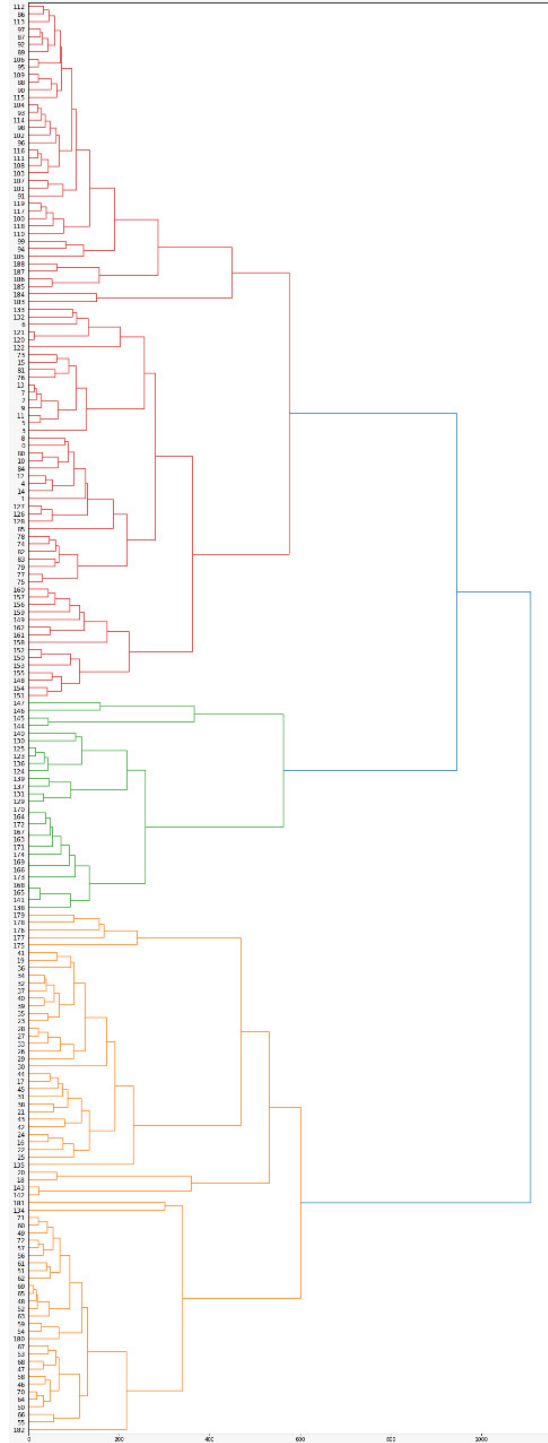


Figure 13: dendrograma con average



A phylogenetic tree showing the relationships between 140 bacterial strains. The tree is rooted on the left and branches out to the right. The strains are color-coded by genus: red for *Bacteroides*, blue for *Bifidobacterium*, green for *Lactobacillus*, and orange for *Streptococcus*. The tree shows a high degree of clustering within each genus, with some inter-genus relationships also visible. The x-axis at the bottom represents genetic distance, with markers at 0.0, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.2, 2.4, 2.6, 2.8, 3.0, 3.2, 3.4, 3.6, 3.8, 4.0, 4.2, 4.4, 4.6, 4.8, 5.0, 5.2, 5.4, 5.6, 5.8, 6.0, 6.2, 6.4, 6.6, 6.8, 7.0, 7.2, 7.4, 7.6, 7.8, 8.0, 8.2, 8.4, 8.6, 8.8, 9.0, 9.2, 9.4, 9.6, 9.8, 10.0, 10.2, 10.4, 10.6, 10.8, 11.0, 11.2, 11.4, 11.6, 11.8, 12.0, 12.2, 12.4, 12.6, 12.8, 13.0, 13.2, 13.4, 13.6, 13.8, 14.0, 14.2, 14.4, 14.6, 14.8, 15.0, 15.2, 15.4, 15.6, 15.8, 16.0, 16.2, 16.4, 16.6, 16.8, 17.0, 17.2, 17.4, 17.6, 17.8, 18.0, 18.2, 18.4, 18.6, 18.8, 19.0, 19.2, 19.4, 19.6, 19.8, 20.0, 20.2, 20.4, 20.6, 20.8, 21.0, 21.2, 21.4, 21.6, 21.8, 22.0, 22.2, 22.4, 22.6, 22.8, 23.0, 23.2, 23.4, 23.6, 23.8, 24.0, 24.2, 24.4, 24.6, 24.8, 25.0, 25.2, 25.4, 25.6, 25.8, 26.0, 26.2, 26.4, 26.6, 26.8, 27.0, 27.2, 27.4, 27.6, 27.8, 28.0, 28.2, 28.4, 28.6, 28.8, 29.0, 29.2, 29.4, 29.6, 29.8, 30.0, 30.2, 30.4, 30.6, 30.8, 31.0, 31.2, 31.4, 31.6, 31.8, 32.0, 32.2, 32.4, 32.6, 32.8, 33.0, 33.2, 33.4, 33.6, 33.8, 34.0, 34.2, 34.4, 34.6, 34.8, 35.0, 35.2, 35.4, 35.6, 35.8, 36.0, 36.2, 36.4, 36.6, 36.8, 37.0, 37.2, 37.4, 37.6, 37.8, 38.0, 38.2, 38.4, 38.6, 38.8, 39.0, 39.2, 39.4, 39.6, 39.8, 40.0, 40.2, 40.4, 40.6, 40.8, 41.0, 41.2, 41.4, 41.6, 41.8, 42.0, 42.2, 42.4, 42.6, 42.8, 43.0, 43.2, 43.4, 43.6, 43.8, 44.0, 44.2, 44.4, 44.6, 44.8, 45.0, 45.2, 45.4, 45.6, 45.8, 46.0, 46.2, 46.4, 46.6, 46.8, 47.0, 47.2, 47.4, 47.6, 47.8, 48.0, 48.2, 48.4, 48.6, 48.8, 49.0, 49.2, 49.4, 49.6, 49.8, 50.0, 50.2, 50.4, 50.6, 50.8, 51.0, 51.2, 51.4, 51.6, 51.8, 52.0, 52.2, 52.4, 52.6, 52.8, 53.0, 53.2, 53.4, 53.6, 53.8, 54.0, 54.2, 54.4, 54.6, 54.8, 55.0, 55.2, 55.4, 55.6, 55.8, 56.0, 56.2, 56.4, 56.6, 56.8, 57.0, 57.2, 57.4, 57.6, 57.8, 58.0, 58.2, 58.4, 58.6, 58.8, 59.0, 59.2, 59.4, 59.6, 59.8, 60.0, 60.2, 60.4, 60.6, 60.8, 61.0, 61.2, 61.4, 61.6, 61.8, 62.0, 62.2, 62.4, 62.6, 62.8, 63.0, 63.2, 63.4, 63.6, 63.8, 64.0, 64.2, 64.4, 64.6, 64.8, 65.0, 65.2, 65.4, 65.6, 65.8, 66.0, 66.2, 66.4, 66.6, 66.8, 67.0, 67.2, 67.4, 67.6, 67.8, 68.0, 68.2, 68.4, 68.6, 68.8, 69.0, 69.2, 69.4, 69.6, 69.8, 70.0, 70.2, 70.4, 70.6, 70.8, 71.0, 71.2, 71.4, 71.6, 71.8, 72.0, 72.2, 72.4, 72.6, 72.8, 73.0, 73.2, 73.4, 73.6, 73.8, 74.0, 74.2, 74.4, 74.6, 74.8, 75.0, 75.2, 75.4, 75.6, 75.8, 76.0, 76.2, 76.4, 76.6, 76.8, 77.0, 77.2, 77.4, 77.6, 77.8, 78.0, 78.2, 78.4, 78.6, 78.8, 79.0, 79.2, 79.4, 79.6, 79.8, 80.0, 80.2, 80.4, 80.6, 80.8, 81.0, 81.2, 81.4, 81.6, 81.8, 82.0, 82.2, 82.4, 82.6, 82.8, 83.0, 83.2, 83.4, 83.6, 83.8, 84.0, 84.2, 84.4, 84.6, 84.8, 85.0, 85.2, 85.4, 85.6, 85.8, 86.0, 86.2, 86.4, 86.6, 86.8, 87.0, 87.2, 87.4, 87.6, 87.8, 88.0, 88.2, 88.4, 88.6, 88.8, 89.0, 89.2, 89.4, 89.6, 89.8, 90.0, 90.2, 90.4, 90.6, 90.8, 91.0, 91.2, 91.4, 91.6, 91.8, 92.0, 92.2, 92.4, 92.6, 92.8, 93.0, 93.2, 93.4, 93.6, 93.8, 94.0, 94.2, 94.4, 94.6, 94.8, 95.0, 95.2, 95.4, 95.6, 95.8, 96.0, 96.2, 96.4, 96.6, 96.8, 97.0, 97.2, 97.4, 97.6, 97.8, 98.0, 98.2, 98.4, 98.6, 98.8, 99.0, 99.2, 99.4, 99.6, 99.8, 100.0, 100.2, 100.4, 100.6, 100.8, 101.0, 101.2, 101.4, 101.6, 101.8, 102.0, 102.2, 102.4, 102.6, 102.8, 103.0, 103.2, 103.4, 103.6, 103.8, 104.0, 104.2, 104.4, 104.6, 104.8, 105.0, 105.2, 105.4, 105.6, 105.8, 106.0, 106.2, 106.4, 106.6, 106.8, 107.0, 107.2, 107.4, 107.6, 107.8, 108.0, 108.2, 108.4, 108.6, 108.8, 109.0, 109.2, 109.4, 109.6, 109.8, 110.0, 110.2, 110.4, 110.6, 110.8, 111.0, 111.2, 111.4, 111.6, 111.8, 112.0, 112.2, 112.4, 112.6, 112.8, 113.0, 113.2, 113.4, 113.6, 113.8, 114.0, 114.2, 114.4, 114.6, 114.8, 115.0, 115.2, 115.4, 115.6, 115.8, 116.0, 116.2, 116.4, 116.6, 116.8, 117.0, 117.2, 117.4, 117.6, 117.8, 118.0, 118.2, 118.4, 118.6, 118.8, 119.0, 119.2, 119.4, 119.6, 119.8, 120.0, 120.2, 120.4, 120.6, 120.8, 121.0, 121.2, 121.4, 121.6, 121.8, 122.0, 122.2, 122.4, 122.6, 122.8, 123.0, 123.2, 123.4, 123.6, 123.8, 124.0, 124.2, 124.4, 124.6, 124.8, 125.0, 125.2, 125.4, 125.6, 125.8, 126.0, 126.2, 126.4, 126.6, 126.8, 127.0, 12

4.4 GMM

Los resultados obtenidos con GMM no reflejaban su funcionamiento correcto. El código se encuentra en el repositorio anexado.

5 Conclusiones

En conclusión, podemos afirmar que el rendimiento de cada algoritmo es diferente y puede variar significativamente según la implementación utilizada. Por ejemplo, AHC es excelente cuando no necesariamente tenemos clusters circulares y no sabemos el número de clusters de antemano. Con este podemos decidir el número de clusters cortando el dendograma, pero con una complejidad cuadrática en la mayoría de casos.

Por otro lado, el DBSCAN es bueno para agrupar puntos estén muy cerca y expendiendo los clusters en dirección a puntos cercanos, lidiando así con diferentes formas de clusters.

Al final, como mencionamos previamente, el rendimiento de cada uno de estos algoritmos está ligado fuertemente al dataset y para que se le decida utilizar.

6 Anexos

- Github: https://github.com/mauriciopinto/proyecto3_IA
- Kaggle: <https://www.kaggle.com/flrotm/clustering>

7 Bibliografía

[1] <https://www.cienciadedatos.net/documentos/py19-pca-python.html>

[2] https://rstudio-pubs-static.s3.amazonaws.com/238819_bc76e21ff51248f7833cf5c90e68197a.html