

Métodos no supervisados

Jordi Casas Roma

1 crédito

xxx/xxxxx/xxx



Índice

| | |
|--|----|
| Introducción | 5 |
| Objetivos | 6 |
| 1. Introducción | 7 |
| 2. Agrupamiento jerárquico | 8 |
| 2.1. Dendrogramas | 9 |
| 2.2. Algoritmos aglomerativos | 9 |
| 2.3. Algoritmos divisivos | 11 |
| 2.4. Ejemplo de aplicación | 12 |
| 2.5. Resumen | 12 |
| 3. El método <i>k-means</i> y derivados | 14 |
| 3.1. <i>k-means</i> | 14 |
| 3.2. Métodos derivados de <i>k-means</i> | 18 |
| 3.3. <i>Fuzzy C-Means</i> | 19 |
| 4. Algoritmo de agrupamiento Canopy | 22 |
| 4.1. El concepto de pre-clustering | 22 |
| 4.2. Funcionamiento del método | 22 |
| 4.3. Ejemplo de aplicación | 24 |
| 4.4. Resumen | 25 |
| 5. HDBScan | 26 |
| Resumen | 27 |
| Glosario | 28 |
| Bibliografía | 29 |

Introducción

El cuarto bloque de este libro se centra en los métodos de aprendizaje no supervisado, es decir, aquellos métodos que utilizan conjuntos de datos no etiquetados. En concreto, se verán algoritmos jerárquicos, el método *k-means* y sus derivados y el concepto y proceso de pre-clustering empleando el algoritmo Canopy.

Objetivos

En los materiales didácticos de este módulo encontraremos las herramientas indispensables para asimilar los siguientes objetivos:

1. XXX

1. Introducción

El cuarto bloque de este libro se centra en los métodos de aprendizaje no supervisado, es decir, aquellos métodos que utilizan conjuntos de datos no etiquetados. En concreto, se verán algoritmos jerárquicos, el método *k-means* y sus derivados y el concepto y proceso de pre-clustering empleando el algoritmo Canopy.

Recordar implicacions no supervisat

Recordar notació de dades D

2. Agrupamiento jerárquico

El agrupamiento jerárquico (*hierarchical clustering* o HC, en inglés) es un método de clasificación no supervisada que busca construir una jerarquía de particiones o clusters en el conjunto de datos.

Hay dos tipos de algoritmos de agrupamiento jerárquico, que aunque pueden proporcionar resultados similares, tienen enfoques inversos:

- El algoritmo aglomerativo (*Agglomerative Hierarchical Clustering, AHC*) parte de una fragmentación completa de los datos, es decir, cada dato tiene su propio grupo, y fusiona los grupos progresivamente hasta que todos los datos pertenecen a un único grupo. Aplica una estrategia *bottom-up*. En este caso hablaremos de clustering o agrupamiento.
- El algoritmo divisivo (*Divisive Hierarchical Clustering, DHC*) parte de un único grupo que contiene todos los datos y lo va dividiendo de forma recursiva hasta obtener un grupo para cada dato. Es decir, aplica una estrategia *top-down*, procediendo de forma inversa a los algoritmos aglomerativos. En este caso hablaremos de segmentación.

Conceptualmente ambos tipos de algoritmos, los aglomerativos y los divisivos, son equivalentes. Sin embargo, los algoritmos aglomerativos son de más fácil construcción, dado que sólo existe un modo de unir dos conjuntos, mientras que existen múltiples formas de separar un conjunto de más de dos elementos.

A través de los algoritmos aglomerativos, es posible construir recomendadores basados en modelo, de modo que para producir una recomendación, sólo es necesario asignar una instancia nueva a uno de los grupos generados por el modelo.

Merece la pena observar que para llevar a cabo la operación de recomendación, tan solo será necesario almacenar la descripción de los grupos generados en el modelo, y no en todo el juego de datos entero, suponiendo así un ahorro notable en recursos.

El principal inconveniente de este tipo de algoritmos es que no son capaces de determinar, por sí mismos, el número idóneo de grupos a generar, sino que es necesario fijarlos de antemano, o bien definir algún criterio de parada en el proceso de construcción jerárquica.

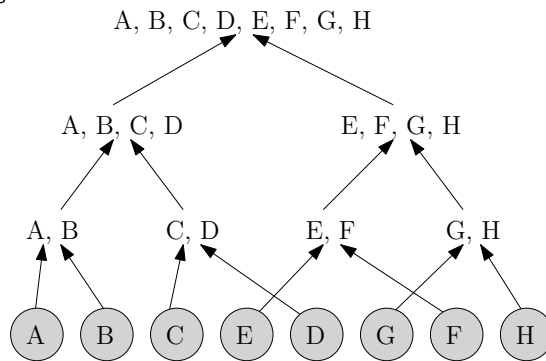
2.1. Dendrogramas

El dendrograma es un diagrama que muestra las agrupaciones (divisiones) sucesivas que genera un algoritmo jerárquico.

Sin duda, es una forma muy intuitiva de representar el proceso de construcción de los grupos. Sin embargo, un punto débil es no ofrecer información sobre la distancia entre los distintos objetos, aunque se puede utilizar el tamaño de las flechas o trabajar con tonalidades de un mismo color para añadir este tipo de información.

En la figura 1 vemos un ejemplo de dendrograma, en el que tenemos una colección de 8 letras y podemos ver como se agrupan progresivamente (algoritmo aglomerativo, estrategia *bottom-up*) o como se dividen recursivamente (algoritmo divisivo, estrategia *top-down*).

Figura 1. Dendrograma aglomerativo. Si las flechas fueran en sentido contrario, tendríamos un dendrograma divisivo



El principal interés de los algoritmos jerárquicos es que generan diferentes niveles de agrupamiento, lo que proporciona una información igual o más útil, en algunos casos, que la obtención del agrupamiento definitivo. Es importante destacar que los algoritmos jerárquicos son **voraces** (*greedy*), lo que significa que buscan la mejor decisión local en cada momento, sin asegurar una solución global óptima.

2.2. Algoritmos aglomerativos

En el Algoritmo 1 podemos ver el detalle del pseudo-código del algoritmo aglomerativo. Como se puede observar, el principal punto de discusión es la medida de la distancia entre dos grupos. Hay dos criterios que afectan ésta medida:

- 1) ¿Cómo se calcula la distancia entre dos puntos? Para éste cálculo podemos utilizar las **métricas de distancia** vistas en la Sección ??, como por ejemplo la distancia euclidiana, distancia de Gauss o distancia de Mahalanobis.
- 2) ¿Como determinamos la distancia mínima entre dos grupos que contienen más de un punto? Cuando tenemos grupos formados por varios puntos, podemos determinar

la distancia entre los dos grupos según diferentes expresiones, que reciben el nombre de **criterios de enlace**.

Para construir un dendograma aglomerativo deberemos, inicialmente establecer con qué métrica de distancia desearemos trabajar y qué criterio de enlace de grupos o segmentos utilizaremos.

El siguiente paso será considerar cada instancia del juego de datos como un grupo o segmento en sí mismo y a partir de aquí empezaremos a calcular distancias entre grupos. En este punto entraremos en un proceso iterativo en el que en cada repetición fusionaremos los grupos más cercanos.

Algoritmo 1 Pseudocódigo del agrupamiento jerárquico aglomerativo (AHC)

Require: D (conjunto de datos individuales)

```

1: Inicializar  $\ell = 0$ 
2: Inicializar el conjunto de clusters  $C_\ell = D$ 
3: while  $|C_\ell| > 1$  do
4:   Encontrar  $c_i, c_j \in C_\ell$  tal que la distancia  $d(c_i, c_j)$  sea mínima
5:   Crear  $c_* = c_i \cup c_j$  como padre de los clusters  $c_i$  i  $c_j$ 
6:   Actualizar  $C_{\ell+1} = C_\ell - \{c_i, c_j\} \cup c_*$ 
7:    $\ell = \ell + 1$ 
8: end while
   return  $C$ 

```

Algunos de los criterios de enlace más utilizados para medir la distancia entre dos grupos A y B son los siguientes:

1) Enlace simple (*simple linkage*): Tomaremos como criterio la distancia mínima entre elementos de los grupos:

$$d(A, B) = \min\{d(a_i, b_j) \mid a_i \in A, b_j \in B\} \quad (1)$$

Puede ser apropiado para encontrar grupos de forma no elíptica. Sin embargo, es muy sensible al ruido en los datos y puede llegar a provocar el efecto cadena. Éste consiste en el hecho de que puede llegar a forzar la unión de dos grupos que a priori deberían permanecer bien diferenciados, por el hecho de que estos compartan algún elemento muy próximo.

2) Enlace completo (*Complete linkage*): Tomaremos como criterio la distancia máxima entre elementos de los grupos:

$$d(A, B) = \max\{d(a_i, b_j) \mid a_i \in A, b_j \in B\} \quad (2)$$

No produce el efecto cadena, pero es sensible a los valores *outliers*. Sin embargo, suele dar mejores resultados que el criterio simple.

3) Enlace medio (*average linkage*): Tomaremos como criterio la distancia media entre elementos de los grupos:

$$d(A, B) = \frac{1}{|A||B|} \sum_{a_i \in A} \sum_{b_j \in B} d(a_i, b_j) \quad (3)$$

Se trata de un criterio que trata de mitigar los inconvenientes de los dos anteriores, sin acabar de resolverlos por completo.

4) Enlace centroide (*centroid linkage*): La distancia entre dos grupos será la distancia entre sus dos centroides. Presenta la ventaja de que su coste computacional es muy inferior al de los criterios anteriores, de modo que está indicado para juegos de datos de gran volumen.

El enlace simple asume, de forma implícita, que la medida de similitud debería de ser transitiva, lo cual es sólo aplicable en algunos dominios. En general, es preferible escoger el enlace completo, que promueve la formación de clusters compactos y homogéneos, aunque pierde la propiedad de monotonicidad. El enlace medio puede representar un buen compromiso entre ambos.

2.3. Algoritmos divisivos

Aunque el agrupamiento aglomerativo representa la forma más natural de construir un dendograma, su coste computacional es alto. El agrupamiento divisivo suele ser más eficiente a nivel computacional, principalmente porque no se suele llegar al nivel mínimo del dendograma (es decir, una clase por dato), dado que no aporta información útil. Se define un criterio de parada que determinará el nivel o profundidad máximo del árbol.

Algoritmo 2 Pseudocódigo del agrupamiento jerárquico divisivo (DHC)

Require: D (conjunto de datos individuales)

Crear cluster inicial $C = \{D\}$ y marcarlo como “abierto”

while (existe un cluster abierto) **do**

 Encontrar el cluster abierto $c^d \in C$

if (c^d no cumple la condición de parada) **then**

 Dividir $c^d = c_1^d \cup c_2^d \cup \dots \cup c_m^d$

for $i = 1, 2, \dots, m$ **do**

 Crear el cluster c_i^d como descendiente de c^d

end for

 Marcar c^d como “cerrado”

else

 Marcar c^d como “cerrado”

end if

end while

El pseudocódigo del algoritmo divisivo se puede ver en el Algoritmo 2. Es importante

destacar los dos criterios principales que afectan el diseño de este tipo de algoritmos: el criterio de parada y el criterio de división de particiones.

El criterio último de parada es cuando una partición tiene un único elemento. En este caso, el algoritmo descendería hasta el nivel de granularidad máximo. Pero generalmente este nivel de detalle no es necesario y añade mucha complejidad de cálculo al algoritmo. Existen tres condiciones de parada básicas, basadas en definir:

- La profundidad máxima del árbol resultante.
- El número mínimo de instancias que deberán de tener las particiones. Las particiones con cardinalidad inferior ya no serán procesadas otra vez.
- Un valor máximo para una medida de distancia o disimilitud. Las particiones que no superen este umbral, no serán divididas otra vez, ya que presentan un nivel de cohesión suficiente.

Referente a los criterios de división, es habitual utilizar algoritmos particionales (no jerárquicos) para crear las divisiones de cada partición. En algunos casos, como por ejemplo en el *k-means*, en el que debemos indicar el número de particiones a realizar, puede ser incluso una ventaja, ya que habilita la construcción de árboles del tamaño deseado (por ejemplo, árboles binarios).

2.4. Ejemplo de aplicación

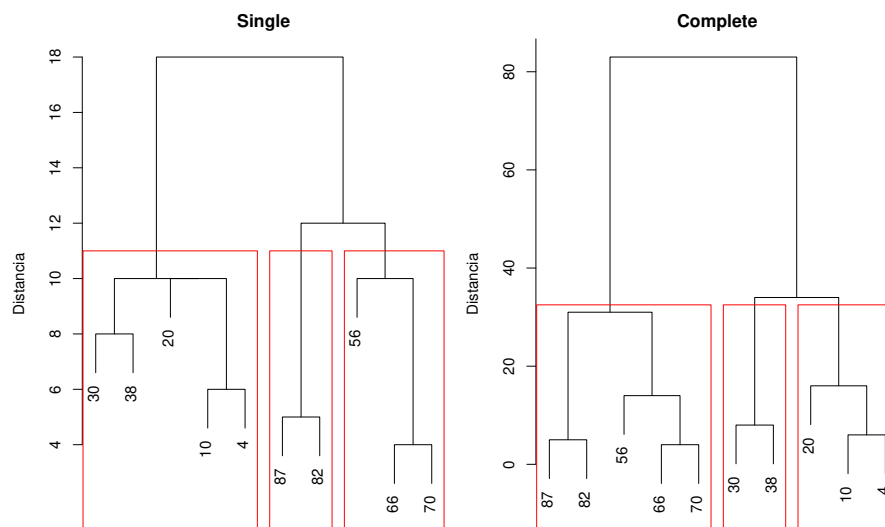
Supongamos que deseamos clasificar un conjunto de 10 números aleatorios entre el 0 y el 100, por ejemplo el conjunto {10, 4, 20, 30, 38, 87, 82, 56, 66, 70}. La Figura 2 muestra los dendogramas generados por el algoritmo aglomerativo utilizando la distancia euclídea y los criterios de enlace simple y completo. El eje vertical nos muestra la distancia entre los distintos números, dando información sobre el grado de similitud entre ellos. Como ya hemos comentado, una de las principales ventajas de los algoritmos jerárquicos es que permiten diferentes niveles de agrupación o granularidad. Por ejemplo, en la Figura 2 podemos escoger un conjunto de tres clusters, lo que genera los grupos marcados (marco rojo).

2.5. Resumen

Los métodos de agrupamiento jerárquico permiten, no sólo descubrir patrones o grupos de similitud en los datos, sino también organizar los datos para un mejor entendimiento. Además, permiten posponer la decisión del nivel de granularidad deseado, permitiendo modificarlo *a posteriori* según los resultados obtenidos.

El agrupamiento jerárquico se suele utilizar como una forma de modelado descriptivo

Figura 2. Ejemplo de dendrogramas aglomerativos utilizando la distancia euclídea y los criterios de enlace simple y completo



más que predictivo. Aunque esta visión es adecuada para la mayoría de aplicaciones, es importante remarcar que el agrupamiento jerárquico también puede ser empleado como modelo predictivo, asignando las nuevas instancias de datos a las hojas o nodos del árbol de agrupación según la métrica de distancia empleada en la construcción del modelo.

Uno de sus mayores problemas es el coste computacional, aunque éste puede ser reducido fijando una profundidad máxima del árbol (en los algoritmos divisivos) o partiendo de una agrupación inicial de instancias (en el caso de los algoritmos aglomerativos).

3. El método *k-means* y derivados

Iniciaremos este capítulo revisando uno de los métodos de clasificación no supervisada más conocidos y empleados, el algoritmo *k-means*. Continuaremos, en la Sección 3.2, presentando algunos de los métodos derivados del algoritmo *k-means* más importantes. Y finalizaremos este capítulo con la versión difusa del algoritmo *k-means*, conocido como *fuzzy C-means*, en la Sección 3.3.

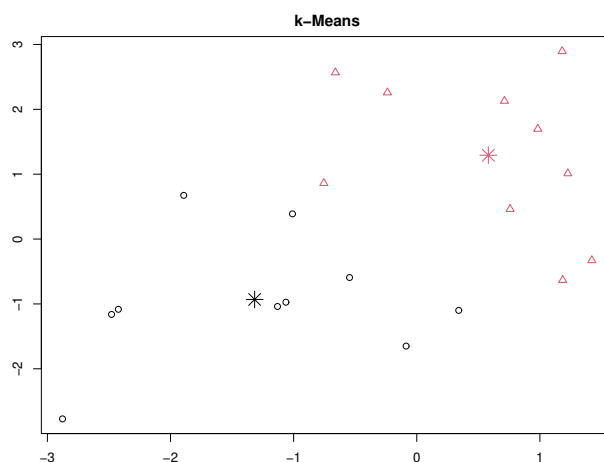
3.1. *k-means*

En esta sección veremos un algoritmo de clasificación no supervisada o segmentación de tipo particional, es decir, que genera una estructura plana. El algoritmo está basado en la partición de un conjunto de n observaciones en k grupos en el que cada observación pertenece al grupo cuya distancia es menor.

Algunas consideraciones importantes respecto al algoritmo *k-means* (o *k-medias*) son:

- El número de grupos o clusters, denotado como k , debe definirse antes de ejecutar el algoritmo.
- Cada grupo o cluster está definido por un punto, generalmente identificado como el centro y llamado semilla o *centroide* del cluster. Aunque puede recibir otros nombres en implementaciones concretas del algoritmo.

Figura 3. Ejemplo de visualización de los resultados de *k-Means*



A grandes rasgos, el algoritmo funciona en dos fases principales, tal y como podemos ver en el pseudocódigo del Algoritmo 3, que corresponden a las dos tareas siguientes:

- 1) En primer lugar, la fase de inicialización, identifica k puntos como *centroides* iniciales. Aunque no es necesario que sean puntos del conjunto de datos, si es importante que sean puntos dentro del mismo intervalo.
- 2) La segunda fase es iterativa, y consiste en
 - a) asignar a cada *centroide* los puntos del conjunto de datos más próximos, formando k grupos disjuntos,
 - b) y a continuación, recalculan los *centroides* en base a los puntos que forman parte de su grupo o partición.

Estos dos pasos se repiten hasta que se satisface la condición de parada.

La Figura 3 muestra el resultado de ejecutar el algoritmo sobre un conjunto de veinte puntos en un espacio bidimensional. Las cruces indican los *centroides* de cada partición.

Algoritmo 3 Pseudocódigo del algoritmo *k-means*

Require: D (conjunto de datos) y k (número de clusters)

```

Inicializar los clusters vacíos  $P = \{p_1, p_2, \dots, p_k\}$ 
Seleccionar los centroides iniciales  $\zeta_1, \zeta_2, \dots, \zeta_k$ 
while (Condición de parada no satisfecha) do
  for all  $d_i \in D$  do
    Asignar la instancia  $d_i$  al cluster  $p_j \mid d(d_i, \zeta_j)$  es mínima
  end for
  for all  $p_i \in P$  do
    Recalcular  $\zeta_i$  según los valores  $d_i \in p_i$ 
  end for
end while
return El conjunto de clusters  $P$ 

```

De los pasos anteriores se desprenden cuatro criterios importantes para definir el comportamiento del algoritmo:

- ¿Cómo podemos inicializar los *centroides*?
- ¿Cómo podemos calcular la distancia entre cada punto d_i y los *centroides*?
- ¿Cómo podemos recalculan los *centroides* a partir de las instancias que forman su grupo o partición?
- ¿Cómo podemos establecer la condición de parada del algoritmo?

3.1.1. Inicialiación de los *centroides*

Un primer enfoque, simple y ampliamente usado, consiste en inicializar los *centroides* con k puntos aleatorios del conjunto de datos. Es decir, $\zeta_i = rand(d_j) \mid d_j \in D, i = \{1, \dots, k\}$.

Una de las principales consecuencias derivadas de este tipo de inicialización de *centroides* es que implica que el proceso de clustering será **no determinista**. Es decir, diferentes ejecuciones del algoritmo pueden conducir a diferentes modelos y, lógicamente, a diferentes niveles de calidad del resultado. Esta característica, que en principio podría parecer un inconveniente importante, puede ser superada gracias a la simplicidad y efectividad del método, que permite ejecutar el algoritmo múltiples veces con distintos *centroides* y escoger el mejor resultado.

El algoritmo *k-means* utiliza la inicialización de *centroides* aleatoria, pero existen otras aproximaciones más complejas que son utilizadas por métodos derivados del *k-means*. Por ejemplo, se pueden seleccionar los *centroides* iniciales a partir de la distribución de probabilidades de las instancias de D , intentando cubrir todo el rango de valores de los datos, especialmente las zonas con mayor concentración de instancias.

3.1.2. Cálculo de la distancia

El algoritmo *k-means*, generalmente, utiliza la distancia euclídea. Aún así, es posible utilizar algún otro tipo de métrica de similitud (consultar las métricas en la Sección ??).

3.1.3. Recálculo de los *centroides*

El valor de cada *centroide* es calculado como la media (*mean*) de todos los puntos que pertenecen a este segmento (de aquí el nombre del algoritmo, *k-means*). Por lo tanto, este algoritmo sólo es aplicable a atributos continuos. En caso de tener atributos no continuos en el conjunto de datos, deberemos aplicar una transformación previa.

3.1.4. Criterio de parada

La **convergencia** del algoritmo es la condición de parada natural y última. Ésta se alcanza cuando no hay cambios en el recálculo de los *centroides* durante una iteración completa, provocando que no haya alteraciones en la distribución de las instancias en las distintas particiones o clusters. Es decir, se llega a una situación de estabilidad en la distribución de las instancias. Esta condición se puede garantizar después de un número finito de iteraciones, dependiendo de la métrica empleada en el cálculo de la distancia y el recálculo de los *centroides*.

En la práctica, no es necesario que el método converja, y generalmente es suficiente cuando se tiene pequeños cambios en la pertenencia de las instancias en los distintos grupos o particiones, es decir, cuando estamos próximos a una situación de convergencia.

3.1.5. Criterios para seleccionar un valor k

Una de las principales características de *k-means* es que se le deben indicar el número de particiones o clusters a identificar. Aunque esto le permite mantener un nivel de simplicidad y eficiencia elevados, puede ser considerado un inconveniente importante. Identificar el número correcto de particiones (k) no es una tarea elemental, e incluso en muchas ocasiones no se puede extraer directamente del conocimiento del dominio. Aún así, en muchos casos el rango de valores de k a considerar no es muy grande, con lo cual se pueden probar diferentes valores y seleccionar el que proporcione mejores resultados.

Un criterio es minimizar la suma de residuos cuadrados (*residual sum of squares*, RSS), es decir, la suma de las distancias de cualquier vector o instancia a su *centroide* más cercano. Este criterio busca la creación de segmentos lo más compactos posibles.

Otro criterio podría ser el de maximizar la suma de distancias entre segmentos, por ejemplo entre sus centros. En este caso estaríamos priorizando tener segmentos los más alejados posible entre sí, es decir, tener segmentos lo más diferenciados posible.

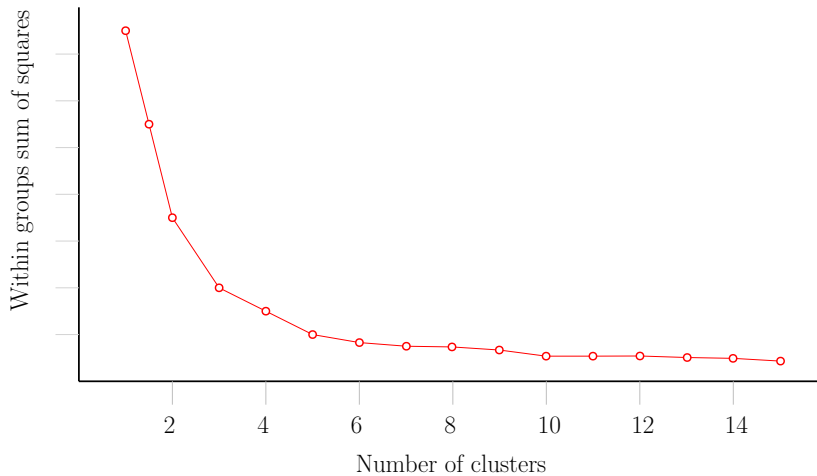
Como sucede en muchos problemas de maximización o minimización, es fácil intuir que el algoritmo no siempre alcanzará un óptimo global, puede darse el caso de que antes encuentre un óptimo local. Además, también es posible que el algoritmo no sea capaz de encontrar ningún óptimo, bien porque simplemente el juego de datos no presente ninguna estructura de segmentos, bien porque no se haya escogido correctamente el número k de segmentos a construir.

Los criterios anteriores (minimización de distancias intra grupo o maximización de distancias inter grupo) pueden usarse para establecer un valor adecuado para el parámetro k . Valores k para los que ya no se consiguen mejoras significativas en la homogeneidad interna de los segmentos o la heterogeneidad entre segmentos distintos, deberían descartarse.

Por ejemplo, en la Figura 4 podemos ver como evoluciona la métrica de calidad respecto al número de particiones creadas. En concreto, observamos cómo a partir de 5 segmentos, la mejora que se produce en la distancia interna de los segmentos ya es insignificante. Este hecho debería ser indicativo de que 5 segmentos es un valor adecuado para k .

Otra aproximación para solucionar este problema, consiste en permitir que el número k se modifique durante la ejecución del algoritmo. En este caso, se inicia el algoritmo con un valor proporcionado de k , pero éste se puede modificar (incrementar o disminuir) según ciertos criterios. Por ejemplo:

- Si existen dos particiones con los centros muy juntos, quizás es mejor unir ambas particiones y reducir el número de particiones.

Figura 4. Criterio para seleccionar un valor k 

- Si existe una partición con un nivel de diversidad muy elevado, se puede dividir ésta en dos nuevas particiones, incrementando por tanto el valor de k .

3.1.6. Resumen

El algoritmo *k-means* es, en esencia, un algoritmo de clasificación o segmentación, lo que permite identificar estructuras dentro de un conjunto de datos sin etiquetas o clases. Aún así, también es ampliamente usado para tareas de predicción, donde se utilizan los *centroides* para clasificar nuevas instancias que no estaban en el conjunto de datos original. Una vez construido el modelo, debe preservarse el valor de los *centroides*, que serán imprescindibles para poder realizar el análisis predictivo.

El algoritmo *k-means* es muy usado como primer intento de clasificación en conjuntos de datos continuos. Por el contrario, no es una buena opción en conjuntos de datos categóricos.

3.2. Métodos derivados de *k-means*

El uso de la distancia media para el cálculo de los *centroides* proporciona un modelo simple y eficiente, y además proporciona unos resultados muy fácilmente interpretables. Aún así, tiene sus limitaciones. Por ejemplo, cuando los datos están distribuidos de forma asimétrica y contienen valores extremos (*outliers*) en algunos atributos.

3.2.1. *k-medians*

El algoritmo *k-mediana* (*k-medians*) es una variación del método anterior, donde se sustituye el cálculo basado en el valor medio, por el valor de la mediana. Aunque la obtención del valor de la mediana requiere mayor complejidad computacional, es preferible en determinados contextos.

Al igual en el caso del *k-means*, éste método puede implementarse mediante distintas métricas de similitud, aunque existe el riesgo de perder la garantía de convergencia. Una de las métricas de distancia más empleadas en este algoritmo es la distancia de Manhattan.

3.2.2. *k-medoids*

El método *k-medians* presenta mejores resultados que el *k-means* cuando las distribuciones son asimétricas o existen valores *outliers*. Aún así, *k-medians* no garantiza que los centros sean similares a alguna instancia del conjunto de datos. Esto se debe a que, a menos que los atributos sean independientes, los vectores de atributos medianos pueden no parecerse a ninguna instancia existente del conjunto de datos ni del dominio completo.

El algoritmo *k-medoids* propone el recálculo de los centros a partir de instancias que presentan un valor de disimilitud mínimo respecto a las demás instancias de la partición o cluster. La identificación de estos puntos, conocidos como *medoids*, es computacionalmente más costosa que las aproximaciones vistas anteriormente.

Aunque esta aproximación sea considerablemente más costosa, a nivel de cálculo, que *k-means* o *k-medians*, es preferible en algunos dominios donde puede ser interesante que los puntos de centro sean puntos del conjunto de datos. Además, este método es particularmente robusto al ruido y a los valores *outliers*.

3.3. *Fuzzy C-Means*

El algoritmo *C-Medios difuso* (*Fuzzy C-Means*) puede ser interpretado como una generalización de los algoritmos vistos anteriormente. La diferencia fundamental, desde el punto de vista operacional, es que este algoritmo construye una partición difusa, mientras que los algoritmos vistos anteriormente construyen una partición nítida o disjunta del conjunto de datos.

La idea que subyace en los métodos difusos se basa en la idea de que no todas las instancias tienen la misma relación de pertenencia con el grupo al que pertenecen. Existen puntos que se encuentran muy próximos al centro de una partición, de los cuales podemos decir que están fuertemente asociados con la partición en cuestión. Otros, por el contrario, se encuentran más próximos al borde que separa dos o más particiones, y por lo tanto, podemos decir que su grado de asociación es más débil que en el caso anterior. Resultaría lógico poder decir que estos segundos mantienen cierta relación con múltiples particiones, aunque pueda haber una partición principal. Esta es la idea principal subyacente en los métodos de segmentación o clustering difuso.

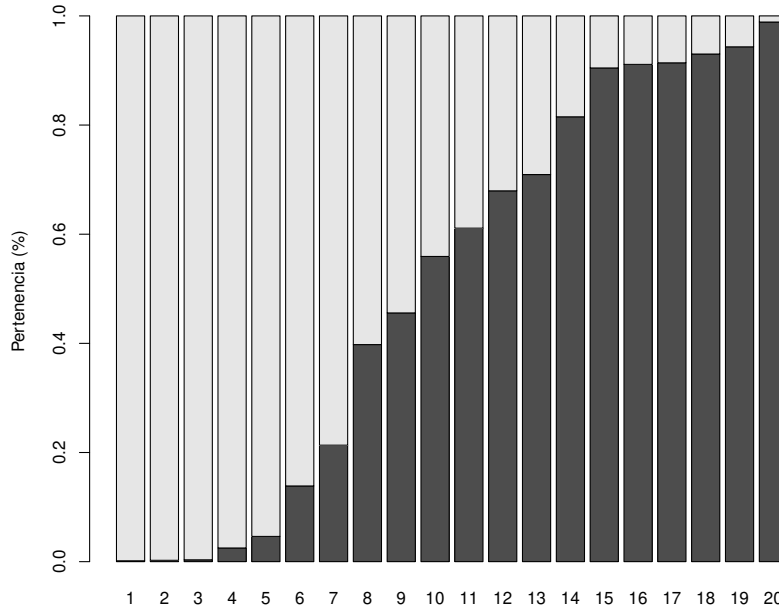
3.3.1. Funcionamiento del método

En este caso, el algoritmo es prácticamente idéntico al Algoritmo 3 visto anteriormente, pero con algunas diferencias que detallamos a continuación:

- Cada instancia d_i tiene asociado un vector de k valores reales que indican el grado de pertenencia de esta instancia a cada uno de los k grupos o particiones. El grado de pertenencia dependerá de la distancia que lo separa con el *centroide* de dicha partición. En general, se normaliza dicho valor para que la suma de todos los valores sea igual a 1. De esta forma, se puede ver el valor como la probabilidad de pertenencia a cada una de las k particiones.
- Al inicio del algoritmo, se asigna el grado de pertenencia de cada instancia a cada grupo aleatoriamente, y después se calculan los *centroides* a partir de estos datos.
- El recálculo de los centroides está ponderado por el grado de pertenencia de cada punto a la partición correspondiente.

La Figura 5 muestra una posible visualización de los resultados de aplicar el algoritmo *Fuzzy C-Means* sobre un conjunto de veinte puntos en el espacio bidimensional. Para cada punto (eje horizontal) se muestra la probabilidad de pertenencia a cada una de las particiones. Se han ordenado los datos según la probabilidad de pertenencia al primer grupo. Es interesante ver que, aunque se pueden diferenciar las dos clases claramente, algunos datos están fuertemente asociados con una única partición, mientras que otros mantienen relaciones de pertinencia un poco más débiles con ambas clases.

Figura 5. Ejemplo de visualización de los resultados de *Fuzzy C-Means*



3.3.2. Resumen

Esta variante, o generalización, presenta las mismas ventajas e inconvenientes que el algoritmo *k-Means*. Por un lado, la simplicidad y eficiencia son sus principales puntos fuertes. Por el otro, el no determinismo y la excesiva dependencia de la métrica utilizada para calcular las distancias son sus posibles limitaciones o inconvenientes.

4. Algoritmo de agrupamiento Canopy

El algoritmo de agrupamiento Canopy (*Canopy clustering algorithm*) (McCallum et al., 2000) es un método no-supervisado de pre-clustering. El pre-clustering se puede entender como el conjunto de tareas de procesamiento previo a un proceso de agrupamiento o clustering, con el objetivo de acelerar el proceso de clustering, especialmente con grandes conjuntos de datos.

4.1. El concepto de pre-clustering

Este algoritmo es utilizado, generalmente, como pre-procesamiento previo a algoritmos de clustering, como por ejemplo el clustering jerárquico aglomerativo (AHC) o el *k-Means*. En conjuntos grandes de datos puede resultar impracticable aplicar directamente los algoritmos de clustering, especialmente en el caso de agrupamiento jerárquico aglomerativo.

La idea que subyace a esta técnica es que podemos reducir drásticamente el número de cálculos que requeridos en el proceso de clustering posterior, introduciendo un proceso previo de generación de grupos superpuestos (*canopies*) a partir de una métrica más sencilla de calcular (*cheapest metric*). De esta forma, sólo calcularemos distancias con la métrica inicial, más estricta y pesada en cálculos, para los puntos que pertenecen al mismo *canopy*.

Podríamos resumirlo diciendo que previamente, mediante una métrica simple, decidimos qué puntos están definitivamente lejos y en consecuencia, para estos puntos alejados ya no valdrá la pena malgastar más cálculos con una métrica más exigente.

4.2. Funcionamiento del método

El método de agrupamiento Canopy divide el proceso de segmentación en dos etapas:

- En una primera etapa, usaremos una métrica sencilla en cálculos, con el objetivo de generar los *canopies* o subgrupos superpuestos de instancias. Cada instancia pueda pertenecer a más de un *canopy* y, a su vez, todas las instancias tienen que pertenecer, al menos, a un *canopy*.
- En una segunda etapa, utilizaremos un método de segmentación tradicional, como por ejemplo el agrupamiento jerárquico aglomerativo (AHC) o el método *k-means*, pero lo haremos con la restricción de no calcular la distancia entre los puntos que

Lectura complementaria

A. McCallum, K. Nigam, L. H. Ungar (2000). "Efficient clustering of high-dimensional data sets with application to reference matching". In Proc. of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, New York, NY, USA, pp. 169-178. <https://doi.org/10.1145/347090.347123>

no pertenecen al mismo *canopy*.

El método, descrito en el Algoritmo 4, parte del conjunto de datos y dos valores límite (*threshold*): T_1 que indica la distancia máxima de la periferia (*the loose distance*) y T_2 que indica la distancia máxima del núcleo (*the tight distance*), donde $T_1 > T_2$. A partir del conjunto inicial de instancias o puntos, se escoge un punto (aleatoriamente) para formar un nuevo *canopy*. A continuación, se calcula la distancia entre este punto y los demás puntos del conjunto de datos, utilizando una métrica más sencilla de calcular – que hemos llamado *cheapest metric*. Se incluyen en el mismo *canopy* todos los puntos que tengan una distancia inferior al *threshold* T_1 . Además, los puntos que tengan una distancia inferior a T_2 , son eliminados del conjunto de datos. De esta forma, se excluye a estos puntos para formar un nuevo *canopy* y ser el centro de éste. Este proceso se repite de forma iterativa hasta que no quedan puntos en el conjunto de datos.

Algoritmo 4 Pseudocódigo del algoritmo Canopy

Require: D (conjunto de datos), T_1, T_2 | $T_1 > T_2$ (distancias)

```

while ( $D \neq \emptyset$ ) do
  Elegir un  $d_i$  aleatoriamente e inicializar nuevo canopy  $\omega_v = \{d_i\}$ 
  Eliminar  $d_i$  del conjunto de datos  $D = D - \{d_i\}$ 
  for all  $d_j \in D$  do
    Calcular la distancia  $\Delta_{i,j} = d(d_i, d_j)$  con la métrica sencilla
    if ( $\Delta_{i,j} < T_1$ ) then
      Añadir  $d_j$  al canopy, i.e.  $\omega_v = \omega_v \cup \{d_j\}$ 
    end if
    if ( $\Delta_{i,j} < T_2$ ) then
      Eliminar  $d_j$  del conjunto de datos  $D = D - \{d_j\}$ 
    end if
  end for
end while
return El conjunto de canopies

```

Para facilitar la comprensión, vamos a situarnos en los dos casos extremos:

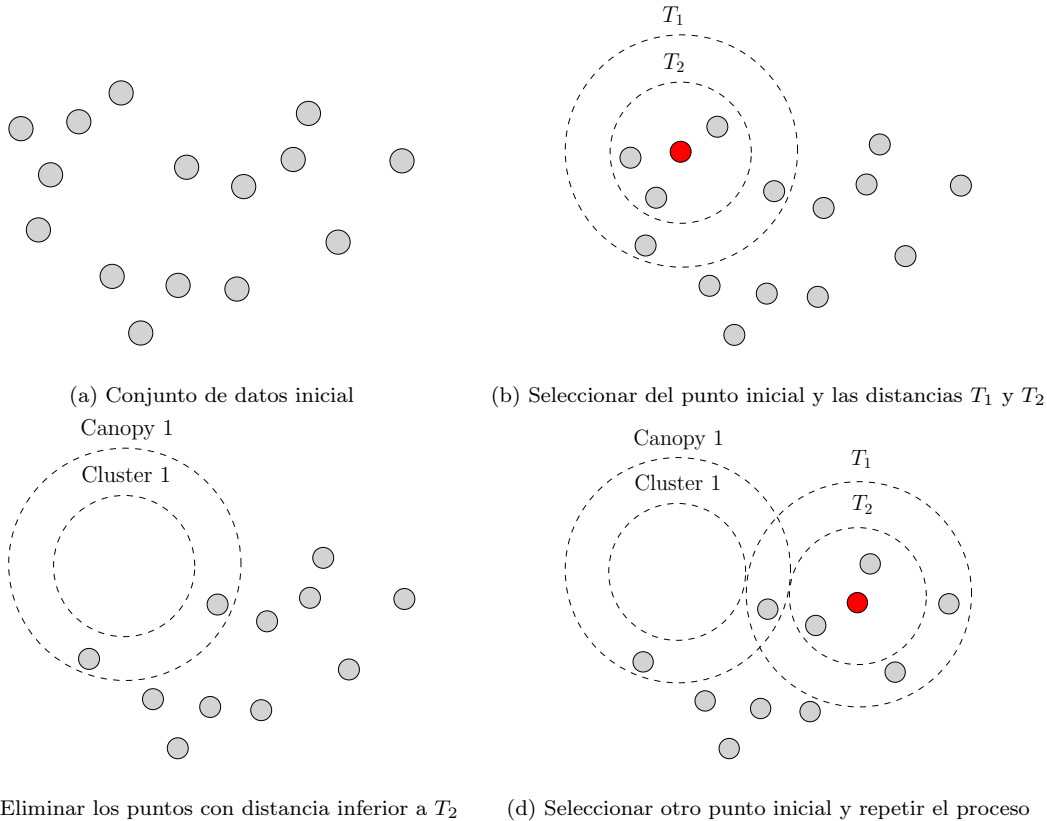
- 1) Supongamos, como consecuencia de la primera etapa, que nuestro universo de puntos cae por completo en un solo *canopy*. Entonces el método de segmentación por *canopies* sería exactamente igual al del método de segmentación tradicional seleccionado, es decir, AHC o *k-means*.
- 2) Supongamos que como resultado de la primera etapa, generamos *canopies* relativamente pequeños y con muy poca superposición. En este caso, al aplicar la técnica tradicional sólo dentro de cada *canopy*, habremos ahorrado un gran número de cálculos.

4.3. Ejemplo de aplicación

Veamos de una forma gráfica cómo funciona el algoritmo:

- 1) En primer lugar, debemos obtener el conjunto de datos D y establecer un valores límite (*threshold*) T_1 y T_2 adecuados para el conjunto de datos con el que estamos trabajado. La Figura 6 (a) muestra un pequeño conjunto de datos, que utilizaremos para ejemplificar.
- 2) A continuación escogemos un punto (aleatoriamente) $d_i \in D$ y lo definimos como el centro del primer *canopy*. Calcularemos la distancia a todos los demás puntos del conjunto de datos D , utilizando la métrica más sencilla de calcular (*cheapest metric*). Por lo tanto, podremos fácilmente definir los puntos que están a distancia menor que T_2 , los que están entre T_2 y T_1 y los que están más allá de T_1 , tal y como muestra la Figura 6 (b).
- 3) Eliminaremos del conjunto de datos D los puntos que están a distancia menor que T_2 , ya que éstos no podrán formar ningún *canopy* nuevo ni formar parte del centro (Figura 6 (c)).
- 4) Repetiremos este proceso forma iterativa hasta que no queden puntos en el conjunto de datos. Es decir, seleccionamos otro punto para forma un nuevo *canopy* y volvemos a aplicar los mismos pasos, manteniendo fijo el valor de T_1 y T_2 (Figura 6 (d)).

Figura 6. Ejemplo de funcionamiento del algoritmo Canopy



4.4. Resumen

El algoritmo de agrupamiento Canopy, o en general los métodos de pre-clustering, presentan la principal ventaja de reducir el número de casos de datos de entrenamiento que se deben comparar en cada paso del método de agrupamiento que se utilice, ya sea *k-means*, agrupamiento jerárquico u otros métodos de agrupamiento no supervisado.

Estos métodos son especialmente indicados en el caso de lidiar con grandes conjuntos de datos, incluyendo conjuntos con millares o millones de instancias y con un número elevado de atributos a considerar.

Aunque se pueda pensar que el uso de este tipo de estrategias de pre-clustering puedan reducir el rendimiento posterior de los métodos de agrupamiento, no se observa un deterioro en los clusters resultantes e incluso, en algunos casos, los resultados mejoran con el proceso de pre-clustering.

5. HDBScan

Resumen

XXX

Glosario

Aprendizaje automático El aprendizaje automático (más conocido por su denominación en inglés, *machine learning*) es el conjunto de técnicas, métodos y algoritmos que permiten a una máquina aprender de manera automática en base a experiencias pasadas.

Aprendizaje no supervisado El aprendizaje no supervisado se basa en el descubrimiento de patrones, características y correlaciones en los datos de entrada, sin intervención externa.

Aprendizaje supervisado El aprendizaje supervisado se basa en el uso de un componente externo, llamado supervisor, que compara los datos obtenidos por el modelo con los datos esperados por éste, y proporciona retroalimentación al modelo para que vaya ajustándose y mejorando las predicciones.

Bibliografía

Y. Bengio, I. Goodfellow, A. Courville (2016). *Deep Learning*. Cambridge, MA: MIT Press.

S. O. Haykin (2009). *Neural Networks and Learning Machines, 3rd Edition*. Pearson.

J. Gironés Roig, J. Casas Roma, J. Minguillón Alfonso, R. Caihuelas Quiles (2017). *Minería de datos: Modelos y algoritmos*. Barcelona: Editorial UOC.