

Análisis y detección de correlaciones en relevamientos transcripcionales
de gran escala

Tesis de Licenciatura en Ciencias Físicas

Facultad de Ciencias Exactas y Naturales

Universidad de Buenos Aires

Andrés Rabinovich

Marzo 2016

Hoja a completar por los jurados

Resumen

Abstract

Dedicatoria

Índice general

1. Métodos de agrupamiento de datos	11
1.1. Similaridad, distancia y disimilaridad	12
1.1.1. Medidas de distancia	13
1.1.2. Similaridad semántica	17
1.2. Estrategias de agrupamiento	20
1.3. Agrupamientos no jerárquicos	20
1.3.1. K-means	21
1.3.2. PAM	21
1.4. Agrupamientos jerárquicos	22
1.4.1. Método de Ward	23
1.4.2. Método de enlace único (o single-link en inglés)	23
1.4.3. Método de enlace completo (o complete-link en inglés)	24
1.4.4. Representación de un agrupamiento jerárquico - dendrogramas	24
1.5. Detectando grupos en el agrupamiento jerárquico	24
1.5.1. Corte de árbol estático	25
1.5.2. Corte de árbol dinámico híbrido	26
1.6. Infomap y CNM	28
1.6.1. Matriz de similaridad topológica - TOM y GTOM	29
2. conclusiones y perspectivas	32

Capítulo 1

Métodos de agrupamiento de datos

Un método de agrupamiento de datos o método de “clustering”, es un método de clasificación no supervisado que permite la partición de un conjunto de N objetos en K grupos o clases, de tal forma que los objetos miembro de un grupo sean más similares entre si (en algún sentido a definir) que entre los miembros de otros grupos.

Son métodos no supervisados ya que en un proceso de agrupación no existen clases definidas previamente ni ejemplos de que tipo de relaciones se desea encontrar entre los objetos, por lo que el mismo proceso debe generar las clases iniciales a las cuales asignar los objetos en el proceso de clasificación.

Estas técnicas permiten el descubrimiento o identificación de distribuciones y patrones subyacentes en los datos, posibilitando obtener conclusiones sobre los mismos, lo que las hace una de las herramientas más útiles en procesos de minería de datos y aprendizaje automatizado en campos tan diversos como las ciencias sociales, las ciencias médicas y la ingeniería.

Dependiendo de los criterios utilizados para realizar la partición, un proceso de agrupamiento puede resultar en diferentes particiones. Como ejemplo de esto podemos tomar el conjunto de números $\{-5, -3, -2, 2, 3\}$. Si decidimos agruparlos por su módulo, obtendremos los conjuntos $\{-5\}$, $\{-3, 3\}$, $\{-2, 2\}$, mientras que si decidimos agruparlos por positividad o negatividad, obtendremos los conjuntos $\{-5, -3, -2\}$ y $\{2, 3\}$. También podríamos haber optado por agrupar por paridad, si son o no primos, etc. Como se observa de un ejemplo tan sencillo, es de fundamental importancia la elección de las propiedades de los objetos a partir de las cuales realizar el agrupamiento. **poner las imágenes de las dos particiones posibles del ejemplo**

En el presente trabajo nos interesará agrupar y caracterizar conjuntos de genes de un organismo modelo, la planta *Arabidopsis thaliana*, en base a sus perfiles de expresión génica a lo largo de diversos tratamientos. [20–22]

Discutiremos a continuación diferentes metodologías y criterios de similaridad que pueden ser considerados para ello.

1.1. Similaridad, distancia y disimilaridad

Las distancias y similitudes tienen un rol preponderante en el análisis de agrupamiento de datos y por regla general son conceptos recíprocos.

Una medida de similitud o coeficiente de similitud se utiliza para indicar de forma cuantitativa la fuerza de la relación entre dos objetos del conjunto. Los $i = 1, 2, \dots, N$ objetos de un conjunto E pueden ser definidos en términos de las coordenadas \vec{X}_i de sus puntos representativos en un espacio $d - dimensional$. Sean $\vec{x} = \{x_0, x_1, \dots, x_d\}$ e $\vec{y} = \{y_0, y_1, \dots, y_d\}$ dos puntos $d - dimensionales$. Entonces, el coeficiente de similitud entre ambos será una función de sus atributos:

$$s(\vec{x}, \vec{y}) = s((x_0, x_1, \dots, x_d), (y_0, y_1, \dots, y_d)) \quad (1.1)$$

con s una función simétrica, es decir, $s(\vec{x}, \vec{y}) = s(\vec{y}, \vec{x})$. Cuanto mayor es el coeficiente de similitud, mayor es la similitud entre ambos.

Por otro lado, las medidas de disimilitud o de distancia se comportan de forma inversa, a mayor distancia o disimilitud, más diferentes son dos puntos. Una métrica de distancia es una función $d \in R$ definida sobre un conjunto E que cumple las siguientes propiedades:

1. No-negatividad: $d(\vec{x}, \vec{y}) \geq 0$
2. Reflexividad: $d(\vec{x}, \vec{y}) = 0 \iff \vec{x} = \vec{y}$
3. Conmutatividad: $d(\vec{x}, \vec{y}) = d(\vec{y}, \vec{x})$
4. Desigualdad triangular: $d(\vec{x}, \vec{y}) \leq d(\vec{x}, \vec{z}) + d(\vec{z}, \vec{y})$

con $\vec{x}, \vec{y}, \vec{z}$ objetos arbitrarios del conjunto.

Una medida de disimilitud es una métrica si cumple con las propiedades antes enunciadas.

Aunque no pareciera existir una definición formal de métrica de similitud, Chen y colaboradores definen una métrica de similitud como una función s que cumple:

1. $s(\vec{x}, \vec{y}) = s(\vec{y}, \vec{x})$
2. $s(\vec{x}, \vec{x}) \geq 0$
3. $s(\vec{x}, \vec{x}) \geq s(\vec{x}, \vec{y})$
4. $s(\vec{x}, \vec{x}) = s(\vec{y}, \vec{y}) = s(\vec{x}, \vec{y}) \iff x = y$
5. $s(\vec{x}, \vec{y}) + s(\vec{y}, \vec{z}) \leq s(\vec{x}, \vec{z}) + s(\vec{y}, \vec{y})$

Si bien es deseable que una similaridad o disimilaridad sea una métrica, existen muchas medidas de similaridad o disimilaridad que dan excelentes resultados en técnicas de agrupamiento de datos sin ser métricas, es decir, sin que necesariamente cumplan la desigualdad triangular o el ítem 5 de métrica de similaridad. [23]

Finalmente, los objetos del conjunto pueden ser especificados por medio de una “matriz de distancia” de $N \times N$ cuyos elementos d_{ij} indican la disimilaridad entre los puntos i y j . [20–22, 24]

1.1.1. Medidas de distancia

El análisis de datos de expresión génica se basa principalmente en la comparación de perfiles de expresión génica. Para poder comprarlos, se requiere una medida que cuantifique cuan similares o disimilares son los objetos considerados. La elección de una medida de distancia será entonces de fundamental importancia para lograr agrupamientos que tengan sentido en el contexto de los datos analizados. En las subsiguientes secciones se listarán las medidas de distancia más comúnmente utilizadas en el agrupamiento de datos (no necesariamente de datos de perfiles de expresión).

Distancia euclidiana

La distancia euclidiana es probablemente la distancia más utilizada en el contexto de datos numéricos. Para dos puntos \vec{x} e \vec{y} en un espacio d – *dimensional*, la distancia euclidiana se define como:

$$d_{euc}(\vec{x}, \vec{y}) = \left[\sum_{i=1}^d (x_i - y_i)^2 \right]^{\frac{1}{2}} = [(\vec{x} - \vec{y})(\vec{x} - \vec{y})^T]^{\frac{1}{2}} \quad (1.2)$$

con x_i e y_i los valores de la i ésima componente de \vec{x} e \vec{y} respectivamente.

Distancia Manhattan o Taxicab

La distancia Manhattan o taxicab es llamada así por ser la distancia que debería recorrer un taxi en una ciudad para ir de un punto a otro, suponiendo la ciudad como una cuadrícula perfecta. Para dos puntos \vec{x} e \vec{y} en un espacio d – *dimensional*, la distancia Manhattan se define como:

$$d_{man}(\vec{x}, \vec{y}) = \sum_{i=1}^d |(\vec{x} - \vec{y})| \quad (1.3)$$

Distancia máxima

Para dos puntos \vec{x} e \vec{y} en un espacio d – *dimensional*, la distancia máxima se define como:

$$d_{max}(\vec{x}, \vec{y}) = \max_{1 \leq j \leq n} |x_i - y_i| \quad (1.4)$$

Distancia de Minkowsky

Para dos puntos \vec{x} e \vec{y} en un espacio d – *dimensional*, la distancia de Minkowsky se define como:

$$d_{mink}(\vec{x}, \vec{y}) = [\sum_{i=1}^d (x_i - y_i)^r]^{\frac{1}{r}}, r \geq 1 \quad (1.5)$$

r es el orden de la distancia de Minkowsky. Notar que si tomamos $r = 2, 1$, ínf obtenemos la distancia euclidiana, la Manhattan y la máxima, respectivamente.

Coefficiente de correlación de Pearson

Una de las métricas más utilizadas para medir similaridad entre perfiles de expresión, como los que se observan en la figura 1.1c, es el coeficiente de correlación de Pearson [6]. El coeficiente de correlación fue descubierto originalmente por Bravais en 1846, pero fue Pearson quién demostró que este coeficiente era la mejor correlación posible entre dos secuencias de números [24].

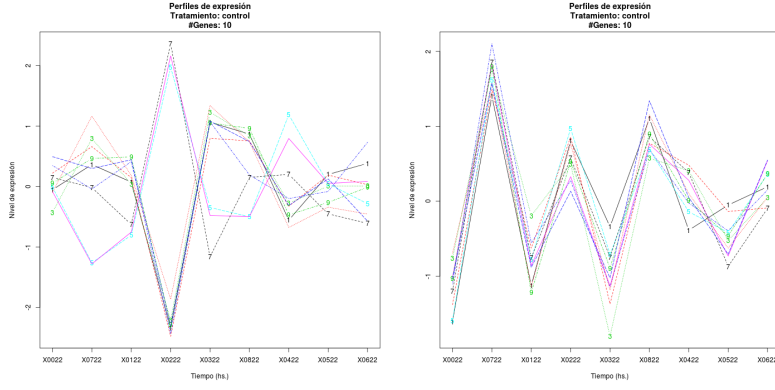
Para dos puntos \vec{x} e \vec{y} en un espacio d – *dimensional*, representando el perfil de expresión de dos genes a lo largo de un dado tratamiento, el CCP se define como:

$$r(\vec{x}, \vec{y}) = \frac{\sum_{i=1}^d (x_i - \bar{x})(y_i - \bar{y})}{[\sum_{i=1}^d (x_i - \bar{x})^2]^{\frac{1}{2}} [\sum_{i=1}^d (y_i - \bar{y})^2]^{\frac{1}{2}}} \quad (1.6)$$

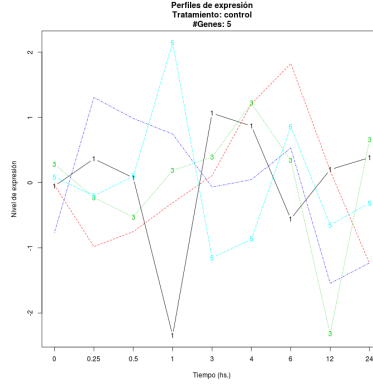
El centrar alrededor de la media permite comparar la forma de ambos perfiles, en lugar de su magnitud.

El coeficiente de correlación r varía entre -1 y $+1$. El caso $r = +1$, llamado *correlación positiva perfecta*, ocurre cuando ambos genes tiene exactamente el mismo perfil, lo que se conoce como **co-regulación hacia arriba positiva (positive up-regulation) ver si este nombre esta bien**, mientras que el caso $r = -1$, llamado *co-regulación hacia abajo negativa perfecta*, ocurre cuando los perfiles son iguales pero opuestos. Un valor del CCP de 0 implica que no se puede inferir una relación entre los perfiles de expresión. La correspondiente medida de distancia puede ser calculada como [25]:

$$d_{ccp}(\vec{x}, \vec{y}) = 1 - r(\vec{x}, \vec{y}) \quad (1.7)$$



(a) Perfiles de expresión para el tratamiento *Control* de 10 genes que están co-regulados. (b) Perfiles de expresión para el tratamiento *Control* de 10 genes que están anti co-regulados.



(c) Perfiles de expresión para el tratamiento *Control* de cinco genes tomados al azar.

Figura 1.1: Distintos grupos de perfiles de expresión

o alternativamente:

$$d_{ccp}(\vec{x}, \vec{y}) = 1 - |r(\vec{x}, \vec{y})| \quad (1.8)$$

En el caso de la distancia definida en 1.8, al tomar el valor absoluto del CCP, genes cuyos perfiles son iguales pero opuestos (están anti co-regulados) pueden encontrarse más cerca en el sentido de d_{ccp} que aquellos que son regulados hacia arriba o abajo pero en distintas magnitudes. Por lo tanto, esta distancia permite encontrar grupos de genes que son co-regulados, sin importar en que sentido (Figura 1.1a) sean co-regulados. En el caso de la distancia definida en 1.7, solamente se consideran cercanos aquellos genes cuyos perfiles sean co-regulados o bien hacia arriba o bien hacia abajo (Figura 1.1b). [26] [24] [6] [20]

En el presente trabajo se utilizará como distancia la definida en 1.7 para encontrar grupos de genes que únicamente se hayan co-regulado o bien hacia arriba o bien hacia abajo [27] discutir o buscar un paper mejor que justifique esto.

1.1.2. Similaridad semántica

La adopción de ontologías provee los medios para comparar aspectos de entidades que de otra forma no podrían ser comparados. Por ejemplo, si dos productos génicos son anotados dentro del mismo esquema, es posible compararlos mediante la comparación de los términos en los cuales están anotados de forma explícita utilizando medidas de similaridad semántica. Se define una medida de similaridad semántica como una función tal que dados dos términos de la ontología o un conjunto de términos en los que dos genes están anotados, la función devuelve un escalar que refleja la cercanía de sentido entre ellos.

Es posible cuantificar la similaridad semántica en una ontología representada por un grafo como GO, mediante diversas estrategias.

Comparación de términos en GO

Existen esencialmente dos formas distintas de comparar términos en GO: comparación a partir de los arcos del grafo y comparación a partir de los nodos y sus propiedades. En este trabajo estaremos interesados únicamente en comparar términos a partir de sus nodos, ya que son los nodos los que contendrán la información biológica en forma de anotaciones génicas.

La comparación a partir de nodos se basa en comparar las propiedades de los términos involucrados, que pueden estar relacionadas con los términos en sí, sus ancestros o sus descendientes. Sea C el conjunto de todos los términos de una ontología GO, con un número total $\#C$ de anotaciones. Un término c_i tendrá $\#c_i$ anotaciones, ya sea directamente o por intermedio de cualquiera de sus hijos. La probabilidad de que un gen tomado al azar, sin otro tipo de información, se encuentre anotado al concepto c_i será entonces $P(c_i) = \frac{\#c_i}{\#C}$, con $P : C \Rightarrow [0 : 1]$.

Se define el contenido de información de c_i como $IC = -\log_2(P(c_i))$, cantidad en el intervalo $(0, -\log_2[\frac{1}{\#C}])$, que indica cuan específico e informativo es un término de la ontología. Para un c_i y c_j tales que $c_i \preceq c_j$, se tiene que $IC(c_i) \geq IC(c_j)$. Cuanto más específico sea un término, es menos probable que un gen dado esté anotado en el mismo, y por lo tanto, su contenido de información es mayor. El nodo raíz de la ontología tiene un contenido de información nulo, ya que es el ancestro de todos los términos de la misma y por lo tanto, saber que un concepto está anotado a la raíz no aporta información.

Si bien el IC puede tener un sesgo, ya que términos en áreas actuales de interés en investigaciones biomédicas van a estar más anotados que otros términos en otras áreas, la utilización del IC sigue teniendo un sentido desde el punto de vista de la probabilidad, porque es mucho más probable (y menos significativo) que dos genes compartan un término frecuentemente usado, que uno no tan frecuente, más allá de si ese término es frecuente porque sea genérico o porque sea un término de interés para la investigación

actual.

Es posible definir una medida entre pares de términos utilizando el IC. El contenido de información puede ser aplicado a los ancestros en común que dos términos poseen, para cuantificar la información que comparten y medir entonces su similaridad semántica. Existen dos formas para ello: tomar el ancestro común más informativo (MICA, por sus siglas en inglés), en donde solo el ancestro común con mayor IC es considerado, o tomar el ancestro disjunto común (DCA por sus siglas en inglés), en la cual todos los ancestros comunes disjuntos (ancestros que no tienen ancestros comunes) son considerados.

Una de las medidas de similaridad semántica más comúnmente utilizadas es la medida de similaridad semántica introducida por Resnik en [16], que consiste en asignar como la medida de similaridad entre dos términos, el contenido de información del primer ancestro en común (el MICA):

$$Sim_{res}(c_i, c_j) = \max_{c \in S(c_i, c_j)} (-\log_2[P(c)]) = IC(MICA[c_i, c_j]) \quad (1.9)$$

Con $S(c_i, c_j)$ el conjunto de ancestros comunes de c_i y c_j .

A modo de ejemplo, tomemos el DAG de la figura 1.2, con 9 términos o conceptos: $C = \{R, c_0, \dots, c_7\}$ y con 5 entidades mapeadas (genes anotados): $g_1 = \{5, 6, 2, 0, r\}$, $g_2 = \{5, 4, 2, 3, 0, r\}$, $g_3 = \{7, 1, r\}$, $g_4 = \{4, 3, 0, r\}$ y $g_5 = \{2, 0, r\}$. Podemos calcular la similaridad semántica de Resnik entre los términos c_4 y c_5 , por ejemplo, sabiendo que $\#C = 5$ y que el ancestro común más informativo de ambos es c_0 , con $\#c_0 = 4$. Se tiene entonces que $Sim_{res}(c_4, c_5) = IC(MICA) = IC(c_0) = -\log_2(\frac{\#c_0}{\#C}) = -\log_2(\frac{4}{5}) = 0,32$. Si quisieramos calcular ahora la similaridad semántica Resnik entre c_5 y c_6 , obtendríamos $Sim_{res}(c_5, c_6) = IC(c_2) = -\log_2(\frac{3}{5}) = 0,73$. Por lo tanto, para Resnik, los conceptos c_5 y c_6 son entre sí, más similares que los conceptos c_4 y c_5 .

Al considerar solo el IC del MICA, la Sim_{res} no tiene en cuenta la especificidad de los términos que compara, es decir, no toma en cuenta la distancia entre los términos y su MICA. Para tomar en cuenta esta distancia, las medidas de Lin [28] y Jiang-Conrath [29] relacionan el IC del MICA con el IC de los términos a comparar:

$$Sim_{lin}(c_i, c_j) = \frac{2 \times IC(MICA[c_i, c_j])}{IC(c_i) + IC(c_j)} \quad (1.10)$$

$$Sim_{JC}(c_i, c_j) = 1 - IC(c_i) + IC(c_j) - 2 \times IC(MICA[c_i, c_j]) \quad (1.11)$$

Eisen1998 Un inconveniente de estas medidas es que se encuentran desplazadas del grafo, es decir, estas medidas son proporcionales a las diferencias entre los IC de los términos y de sus ancestros comunes, independientemente del valor absoluto de IC del ancestro. Una restricción de todas estas medidas es que solo toman en cuenta el MICA, a pesar de que los términos GO pueden tener varios ancestros disjuntos comunes (DCA). Para evitar esta restricción, [17] propuso la aproximación GraSM, que puede ser aplicada a todas las medidas descritas anteriormente, simplemente reemplazando el

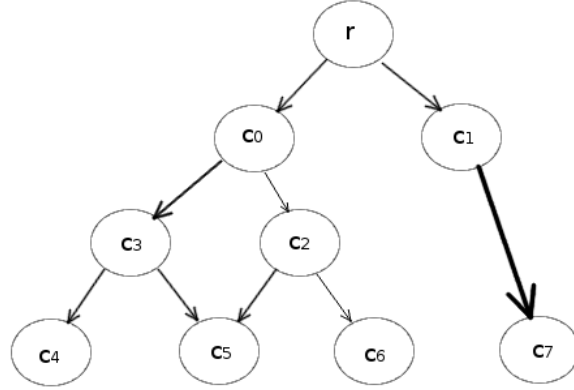


Figura 1.2: DAG con 9 términos o conceptos: $C = \{R, c_0, \dots, c_7\}$ y con 5 entidades mapeadas (genes anotados): $g_1 = \{5, 6, 2, 0, r\}$, $g_2 = \{5, 4, 2, 3, 0, r\}$, $g_3 = \{7, 1, r\}$, $g_4 = \{4, 3, 0, r\}$ y $g_5 = \{2, 0, r\}$

IC del MICA, por el promedio de los IC de los DCA. Existen más de dos docenas de medidas de similitud entre términos GO, y no siempre es claro cuál es el mejor para un dado propósito. Sin embargo, generalmente la elección de una medida por defecto es suficiente [5]. En este trabajo utilizaremos la Sim_{res} , por tratarse de una medida simple y efectiva.

Existen dos estrategias distintas para asignar una similitud semántica entre genes. La primera se basa en medidas globales (groupwise en inglés), que comparan globalmente los conjuntos de términos en los que dos genes están anotados, $GO(g_1)$ y $GO(g_2)$, por ejemplo, contando cuantos términos comparten: $|GO(g_1) \cap GO(g_2)|$. [30]

La segunda estrategia se basa en medidas de a pares (pairwise en inglés), calculando la similitud semántica término a término de cada uno de los conjuntos $GO(g_1)$ y $GO(g_2)$ y luego aplicando sobre esta similitud alguna operación para obtener una medida de similitud entre estos genes.

El primer paso para esto es calcular una matriz de similitud S de $N \times M$ que contenga la similitud de a pares, entre todos los pares de términos de estos conjuntos, con $N = |GO(g_1)|$ y $M = |GO(g_2)|$, utilizando alguna de las medidas de similitud semántica entre términos presentadas anteriormente (Sim_{res} , Sim_{lin} , etc.):

$$S_{ij} = Sim(GO(g_1^i), GO(g_2^j)), \forall i \in \{1, \dots, N\} \forall j \in \{1, \dots, M\} \quad (1.12)$$

Notar que esta matriz puede no ser simétrica.

Cada una de las N filas corresponde a la similitud entre la anotación i –esima del gen 1 y todas las M anotaciones del gen 2 y cada una de las M columnas corresponde a la similitud entre la anotación j –esima del gen 2 y todas las N anotaciones del gen 1.

A partir de S_{ij} es posible definir tres métodos para obtener una medida de similitud

entre genes. El primer método, propuesto en [31], consiste en tomar como similaridad, la máxima similaridad entre todos los términos:

$$Sim_{max}(GO(g_i), GO(g_j)) = \max\{S_{ij}\} \quad (1.13)$$

El segundo método, propuesto en [32], consiste en tomar el valor medio de todos los valores de la matriz S_{ij} :

$$Sim_{med}(GO(g_i), GO(g_j)) = \frac{1}{N.M} \sum_{i,j} S_{ij} \quad (1.14)$$

Finalmente, el tercer método, propuesto en [17], implica tomar el valor medio de los máximos de cada fila, el valor medio de los máximos de cada columna, y quedarse con el máximo de esos dos valores. Este criterio de similaridad se conoce como *rcmax*:

$$Sim_{rcmax}(GO(g_1), GO(g_2)) = \max\left\{\frac{1}{N} \sum_i \max_{1 \leq j \leq M} S_{ij}, \frac{1}{M} \sum_j \max_{1 \leq i \leq N} S_{ij}\right\} \quad (1.15)$$

Como muchos genes están anotados en conceptos muy diversos por participar en procesos biológicos muy distintos, e incluso puede haber genes que no están anotados en ningún concepto, la medida de similaridad Sim_{med} tiende a dar valores más bajos que otros métodos. Por el contrario, la medida Sim_{max} tiende a dar valores más altos, por ser una medida más optimista. En este trabajo utilizaremos el tercer método, Sim_{rcmax} , por ser un compromiso entre ambos casos extremos. [16] [17] [5] [28] [29] [31] [32]

1.2. Estrategias de agrupamiento

En lo que sigue introduciremos las diferentes estrategias de agrupamiento de datos utilizadas en este trabajo, tanto para agrupamiento de perfiles transcripcionales como de armado de comunidades en las redes presentadas anteriormente.

Es posible distinguir dos tipos de agrupamientos, conocidos como agrupamiento duro (hard clustering en inglés), y agrupamiento difuso (fuzzy clustering en inglés). En el primer caso, el de agrupamiento duro, cada objeto del conjunto de datos es asignado a un y solo un grupo, mientras que en el segundo caso, el de agrupamiento difuso, un elemento del conjunto puede pertenecer a varios grupos, con distinta probabilidad. En este trabajo utilizaremos únicamente métodos de agrupamiento duro.

1.3. Agrupamientos no jerárquicos

Además la distinción mencionada más arriba, los métodos de agrupamiento pueden dividirse (entre otros) fundamentalmente entre agrupamientos jerárquicos y agrupa-

mientos no jerárquicos. Las dos estrategias de agrupamientos no jerárquicos que se presentan a continuación fueron utilizadas en el desarrollo de este trabajo.

1.3.1. K-means

K-means es un método usual de agrupamiento no jerárquico en donde cada observación pertenece al grupo con la media más cercana a la observación.

El mismo comienza agrupando los objetos de forma arbitraria en K grupos distintos. El número K puede ser elegido de forma aleatoria o estimado mediante algún otro método de agrupamiento jerárquico pero es siempre fijo. Luego, se calcula un promedio de la posición de todas las observaciones de cada grupo, llamado centroide. A continuación, los objetos individuales son redistribuidos de un grupo a otro dependiendo de que centroide esté más cerca de la observación. Este procedimiento de calcular el centroide de cada cluster y re agrupar los objetos más cercanos a los centroides disponibles se repite de manera iterativa una cantidad fija de veces o hasta la convergencia del método (se considera que el método converge cuando una iteración no modifica la iteración anterior).

Típicamente, se requieren entre 20000 y 100000 iteraciones para la convergencia del método, aunque no hay garantías de que eso ocurra. Más formalmente, sea un conjunto de observaciones $\{\vec{x}_1, \dots, \vec{x}_n\}$, k-means construye una partición de las observaciones en k grupos con $k \leq n$ a fin de minimizar una función de costo, como ser la suma de los cuadrados dentro de cada grupo $G = \{g_1, \dots, g_k\}$:

$$C = \underset{i=1}{\operatorname{argmin}} \sum_{x_j \in g_i}^k ||x_j - \mu_i|| \quad (1.16)$$

Con μ_i el valor medio de los elementos del grupo g_i . La figura 1.3 muestra un conjunto de observaciones y los grupos que se obtienen fijando $k = 2$ y $k = 5$, junto con sus respectivos centroides. Se observa que dependiendo del k utilizado, el algoritmo encuentra particiones con mayor o menor nivel de *resolución*. Volveremos sobre el tema de la resolución más adelante. [33] [34]

1.3.2. PAM

Si bien k-means es uno de los métodos de partición más utilizados ya que es muy eficiente en términos de tiempo computacional, el mismo es muy sensible a observaciones aisladas. Por esta razón, en algunos métodos se reemplazan los centroides, que son puntos no necesariamente pertenecientes al conjunto de observaciones, por medoides, que son los objetos más centrales dentro del grupo (se reemplaza k-means por k-medoids). Esto hace que el método sea insensible a observaciones aisladas.

Particionar alrededor de medoides (Partitioning around medoids en inglés) es uno de los

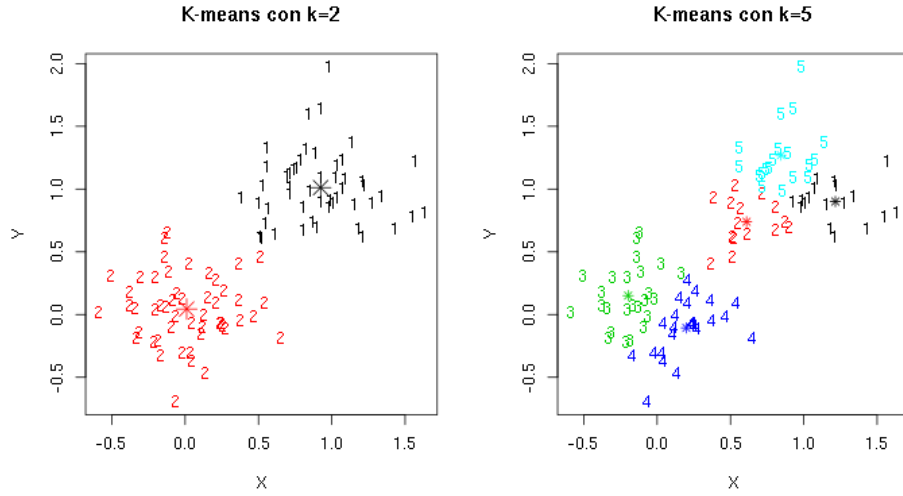


Figura 1.3: Agrupamiento utilizando k-means con $k = 2$ y $k = 5$. **mejorar epigrafe**

métodos más conocidos que hace uso de este concepto, buscando minimizar la función de costo:

$$C = \underset{i=1}{\operatorname{argmin}} \sum_{x_j \in g_i}^k d(x_j, m_i) \quad (1.17)$$

Con m_i el medoide del grupo i y $d(x_j, m_i)$ la distancia entre el objeto x_j del grupo i y el medoide del mismo grupo. [35] [36]

1.4. Agrupamientos jerárquicos

Existen dos acercamientos distintos para realizar un agrupamiento jerárquico: se puede ir “desde abajo hacia arriba”, agrupando grupos más chicos en grupos más grandes, lo que se conoce como agrupamiento aglomerativo, ose puede ir “desde arriba hacia abajo”, dividiendo grupos más grandes en grupos más chicos, lo que se conoce como agrupamiento divisivo. En este trabajo nos interesará únicamente trabajar con agrupamientos aglomerativos.

Un agrupamiento jerárquico aglomerativo comienza con cada objeto en un grupo separado. Luego, se unen los dos grupos más cercanos de acuerdo a algún criterio definido generando un nuevo grupo a partir de ambos. Al nuevo grupo se le asignará una distancia al resto de los grupos de acuerdo a cierto criterio. Esto se repite hasta que solo quede un único grupo.

Es un tipo de procedimiento determinista y voraz (greedy en inglés), ya que realiza las decisiones tomando en cuenta los óptimos locales en cada etapa, esperando obtener con

esto un óptimo global.

Se dice que una partición es más fina (o un refinamiento) de otra partición, si cada grupo de una partición más fina está contenido dentro de un grupo de la partición más gruesa, es decir, cada grupo de la partición más fina es un sub-grupo de un grupo de la partición más gruesa. El agrupamiento jerárquico es un método cuyo resultado es un conjunto de particiones anidadas P_n, P_{n-1}, \dots, P_1 cada vez más gruesas, donde cada nivel más alto une dos grupos de una partición de un nivel más bajo.

Para poder realizar este procedimiento, es necesario definir cuan cercanos son dos grupos:

1.4.1. Método de Ward

Este método busca unir los grupos de una forma tal que se minimice la pérdida de información asociada a cada unión, usualmente cuantificada como el error de la suma de los cuadrados (ESS). Dado un conjunto de puntos C , el ESS asociado a C queda definido por:

$$ESS(C) = \sum_{\vec{x} \in C} (\vec{x} - \mu(C))(\vec{x} - \mu(C))^T \quad (1.18)$$

con $\mu(C) = \frac{1}{|C|} \sum_{\vec{x} \in C} \vec{x}$, el valor medio de C . Suponiendo que una dada partición está separada en k grupos, $\{C_1, C_2, \dots, C_k\}$, entonces se tiene que la pérdida de información de la partición está dada por:

$$ESS = \sum_{i=1}^k ESS(C_i) \quad (1.19)$$

En cada etapa de este método, se prueban todas las uniones de grupos posibles de a pares y se realiza aquella unión que minimiza 1.19.

En el agrupamiento jerárquico, el ESS comienza en cero, ya que cada punto pertenece a un grupo distinto, y crece a medida que se unen grupos. Al ser un algoritmo voraz, la ESS para un dado número de grupos k no será necesariamente la mínima.

1.4.2. Método de enlace único (o single-link en inglés)

Este método es uno de los métodos más simples para agrupamiento jerárquico. El mismo define la distancia entre dos grupos como la mínima distancia entre sus miembros. Sean C_i y C_j dos grupos, entonces la distancia de enlace único se define como:

$$D_{sl}(C_i, C_j) = \min_{\vec{x} \in C_i, \vec{y} \in C_j} d(\vec{x}, \vec{y}) \quad (1.20)$$

con $d(\vec{x}, \vec{y})$ la función de distancia utilizada para calcular la matriz de disimilaridad entre los elementos. El nombre de enlace único hace referencia a que dos grupos están

cerca aunque tengan un único par de puntos cerca. Este método permite el manejo de grupos con formas complejas y es invariante ante transformaciones monótonas (como una transformación logarítmica) [37].

Este algoritmo solamente considera la separación entre elementos, dejando de lado la compacidad o el balance en los grupos.

1.4.3. Método de enlace completo (o complete-link en inglés)

Este método es similar al método de enlace único, ya que toma la distancia entre dos grupos como el máximo de la distancia entre sus puntos:

$$D_{cl}(C_i, C_j) = \max_{\vec{x} \in C_i, \vec{y} \in C_j} d(\vec{x}, \vec{y}) \quad (1.21)$$

con $d(\vec{x}, \vec{y})$ la función de distancia utilizada para calcular la matriz de disimilaridad entre los elementos.

En este trabajo, utilizaremos el método de enlace completo. [38] [20] [37]

1.4.4. Representación de un agrupamiento jerárquico - dendrogramas

Un agrupamiento jerárquico puede representarse como un árbol, llamado dendrograma, que permite una rápida interpretación. En un dendrograma, cada nodo está asociado con una altura h , tal que si A y B son dos nodos del dendrograma, h cumple:

$$h(A) \leq h(B) \Leftrightarrow A \subseteq B \quad (1.22)$$

A modo ilustrativo, la figura 1.4 muestra el agrupamiento jerárquico realizado sobre 10 puntos colocados de forma aleatoria en el plano, agrupados utilizando la distancia euclidiana y mediante los tres métodos vistos anteriormente (Ward, enlace único y enlace completo). De estos gráficos es claro que cada método produce una secuencia diferente de particiones, y dependerá de la aplicación que se requiera, cual de los métodos utilizar.

1.5. Detectando grupos en el agrupamiento jerárquico

El agrupamiento jerárquico organiza los objetos en árboles (dendrogramas) cuyas ramas son los grupos deseados. El proceso de detección de grupos se conoce como corte de árbol, corte de ramas o podado de ramas.

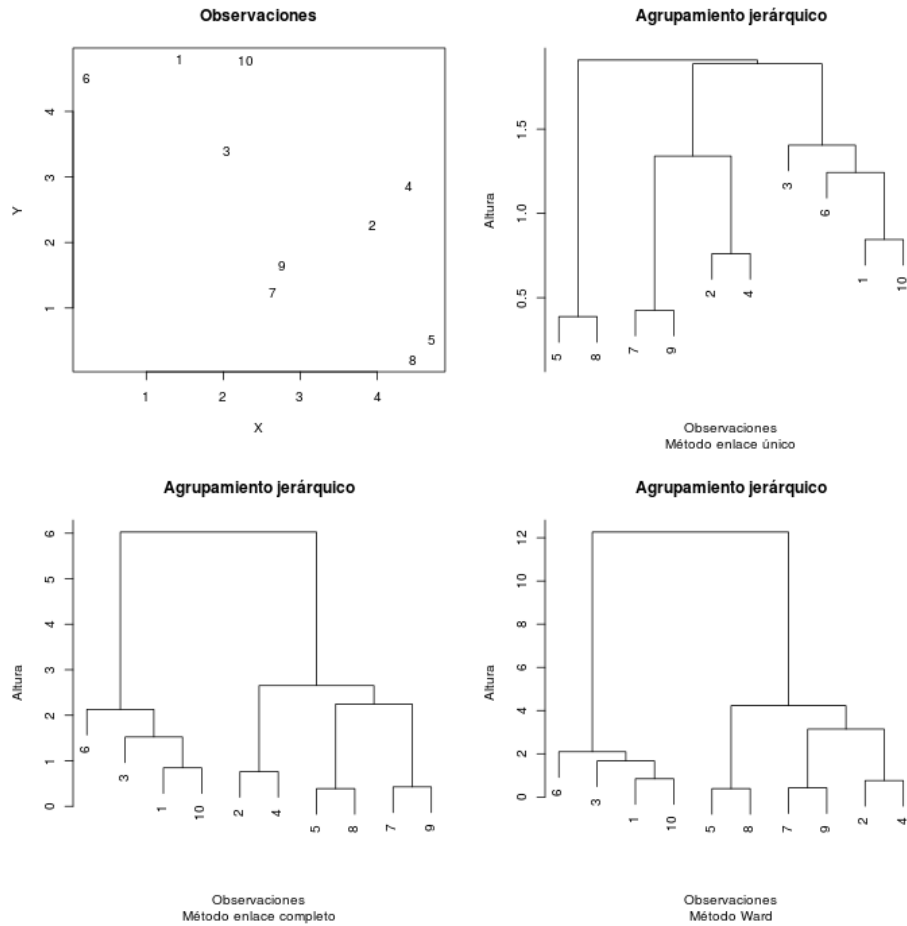


Figura 1.4: Ejemplos de agrupamientos jerárquicos utilizando el mismo conjunto de datos pero distintos métodos de distancia entre grupos. **mejorar epigrafe**

1.5.1. Corte de árbol estático

El método más sencillo de podado es conocido como corte de árbol estático, y funciona definiendo cada rama contigua debajo de una altura fija de corte, como un grupo separado. La cantidad de grupos obtenidos por éste método depende fuertemente de la altura de corte elegida. La figura 1.5 muestra dos alturas de corte posibles y los grupos que se obtienen a partir de cada una de ellas. Al cortar el árbol en $h = 3$, se obtienen dos grupos, el grupo g_1 , que contiene a las observaciones $\{6, 3, 1, 10\}$ y el grupo g_2 que contiene a las observaciones $\{2, 4, 8, 5, 7, 9\}$, mientras que al cortarlo en $h = 2$, se obtienen cuatro grupos, g'_1 con la observación $\{6\}$, g'_2 con las observaciones $\{3, 1, 10\}$, g'_3 con las observaciones $\{2, 4\}$ y g'_4 con las observaciones $\{5, 8, 7, 9\}$.

A partir de un ejemplo tan sencillo es inmediato notar que el problema del agrupamiento es un problema “mal planteado”, es decir, cualquier conjunto de puntos puede ser

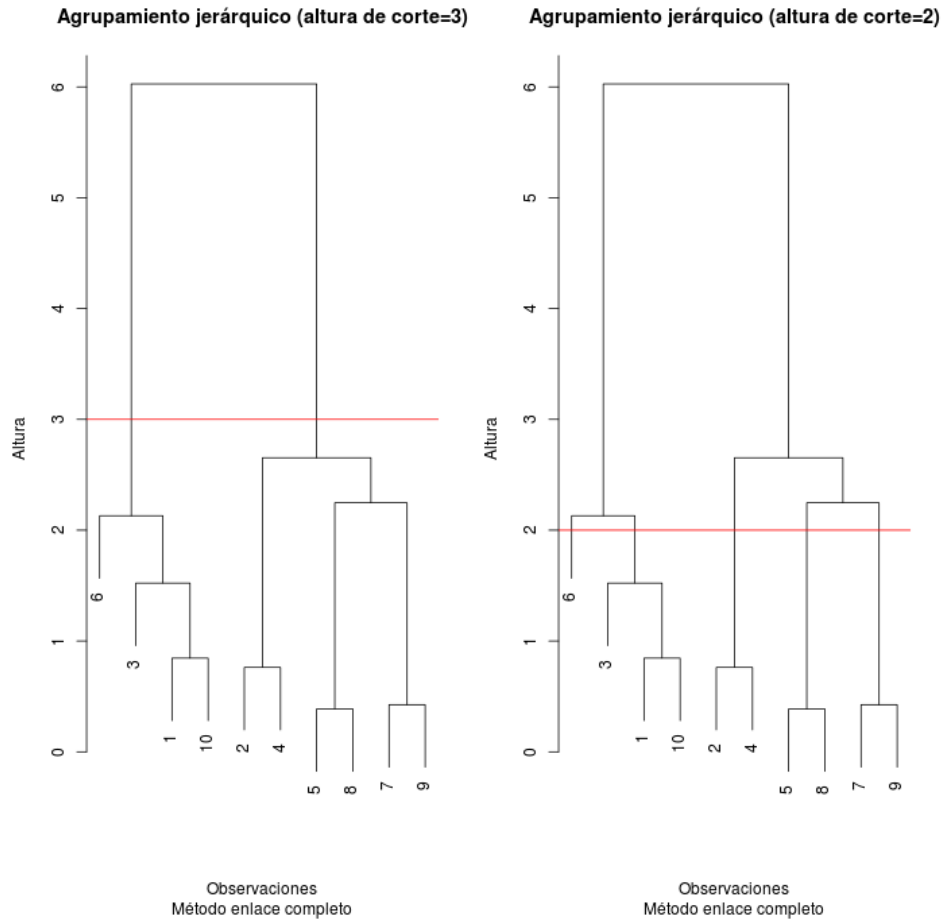


Figura 1.5: Corte de árbol a dos alturas diferentes. **mejorar epigrafe**

agrupado de maneras drásticamente distintas, sin que exista actualmente un criterio para preferir uno u otro agrupamiento. La fuente de ambigüedades a este respecto más importante, es que la forma en que los datos deberían ser agrupados, depende fuertemente de la *resolución* deseada. Lo que parece una única nube de puntos puede resultar ser, al analizar los datos con mayor resolución, una partición compuesta de muchos grupos. Cada tarea deberá encontrar el nivel adecuado de resolución para obtener la cantidad “correcta” de grupos. [22] [39]

1.5.2. Corte de árbol dinámico híbrido

Si bien es posible detectar grupos distintos en el dendrograma a partir de una inspección visual, utilizar una técnica de corte de árbol estático de forma programática no siempre logra identificar adecuadamente los grupos. Este no es un inconveniente

del método de agrupamiento jerárquico, sino que al poseer grupos anidados, un solo corte a una altura prefijada no será capaz de detectarlos todos. El método de corte de árbol dinámico híbrido ataca este problema analizando la forma de las ramas del dendrograma en lugar de una altura absoluta. El mismo construye los grupos de abajo hacia arriba en dos pasos. En el primer paso, se detectan las ramas que satisfacen un criterio específico para ser grupos. Este paso de poda está basado en la información de unión del dendrograma. En el segundo paso, se miden cuán cerca de los grupos detectados en el primer paso están todos los objetos no asignados previamente. Si un objeto está suficientemente cerca de un grupo, es asignado a ese grupo. En este paso, se ignora el dendrograma y se utiliza únicamente la información de disimilaridad. Este paso puede considerarse un método modificado de particionado alrededor de medoides (modified Partitioning Around Medoids o mPAM, en inglés). Por eso el nombre de *híbrido*, al tratarse de una mezcla entre agrupamiento jerárquico y no jerárquico. Los criterios específicos para la detección de grupos se basan en los siguientes cuatro criterios de la forma de las ramas:

1. Un grupo debe tener una cantidad mínima de objetos.
2. Los objetos que están muy lejos del grupo son excluidos del grupo aunque pertenezcan a la misma rama del dendrograma.
3. Cada grupo debe estar separado de su entorno por una brecha o espacio vacío.
4. El núcleo de cada grupo (el conjunto de objetos con menor altura de unión en el grupo) debe estar fuertemente conectado.

O más formalmente, dado un núcleo de un grupo, llamamos d al promedio de las disimilaridades de pares entre objetos del núcleo, es decir, a su dispersión y definimos la brecha g de un grupo como la diferencia entre d y la altura donde el grupo se une al resto del dendrograma y entonces, una rama se considera un grupo si:

1. Tiene al menos N_0 objetos.
2. Todas las alturas de unión son a lo sumo de h_{max} .
3. La brecha g del grupo es mayor que un g_{min} .
4. La dispersión d del núcleo es a lo sumo d_{max} .

Los parámetros N_0 , h_{max} , g_{min} y d_{max} son parámetros ajustables del método. La figura 1.6 muestra un ejemplo de los parámetros utilizados para definir los grupos en el paso 1.

Para el paso 2, de tipo PAM, los objetos no asignados (o aquellos grupos que no cumplan tener al menos N_0 objetos) son asignados al grupo más cercano si la disimilaridad

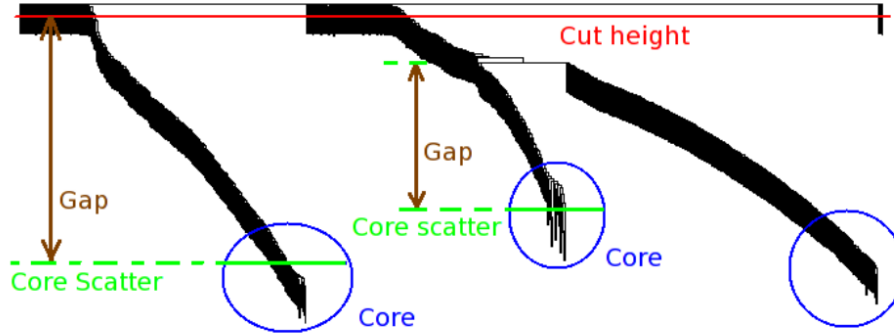


Figura 1.6: Dendrograma simulado con tres ramas con alturas de unión diferentes. La altura de corte corresponde a h_{max} . **poner en castellano y citar correctamente**

correspondiente es más pequeña que una disimilaridad máxima definida previamente, o si es más pequeña que el “radio” del grupo. El “radio” se define como la máxima de las disimilaridades del medoide del grupo al resto de los objetos del mismo.

Es posible controlar la sensibilidad de las divisiones de los grupos mediante el parámetro *deepSplit*, que puede tomar los valores de 1 a 4. Para un *deepSplit* = 1, el método producirá relativamente pocos grupos, de muchos elementos y bien definidos, mientras que para *deepSplit* = 4, el método producirá más grupos pero con una dispersión mayor en el núcleo y separado por brechas más pequeñas.

Para una descripción más detallada del algoritmo, el lector interesado puede referirse a [39], [40].

1.6. Infomap y CNM

Como se desarrolló en la sección ??, las redes son construcciones útiles para esquematizar la organización de las interacciones en distintos tipos de sistemas. Sin embargo, por motivos de visualización, solo se pueden representar pequeños sistemas. Las redes reales son usualmente tan grandes que es necesario representarlas mediante algún mecanismo de granularidad más gruesa, es decir, descomponer a la red en módulos que representen varios nodos y arcos. Este es el objetivo básico de lo que se conoce como *detección de comunidades*.

En este trabajo utilizaremos dos métodos de modularización en redes, Infomap y CNM. El método o algoritmo Infomap hace uso de criterios de optimización basados en teorías de información, donde los módulos se definen de tal forma que la longitud media de la descripción de un proceso de paseo al azar en el grafo sea mínima, mientras que el algoritmo de Clauset-Newman-Moore (CNM), a partir de ciertas heurísticas, busca particiones de la red optimizando directamente una función de calidad Q .

Ambos métodos serán utilizados en este trabajo con el fin de comparar los resultados obtenidos para las comunidades Infomap y CNM con los obtenidos para los métodos de agrupamiento usados. [18] [41]

Distintos métodos darán distintos resultados, dependiendo del conjunto de datos y del objetivo del agrupamiento, por lo que es de vital importancia elegir el método adecuado a la aplicación en cuestión. Las figuras 1.7 y 1.8 muestran ejemplos de conjuntos de datos diversos y de como son agrupados por los distintos métodos.

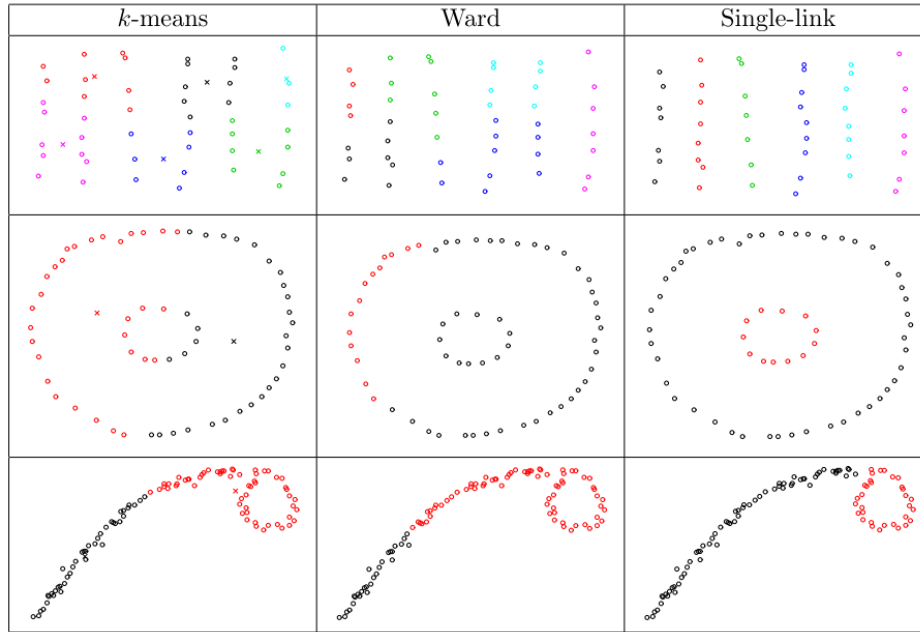


Figura 1.7: Algunos casos para los cuales el método de enlace único (single-link) se comporta “mejor” que los métodos de k -means o de Ward .[citar correctamente](#)

1.6.1. Matriz de similaridad topológica - TOM y GTOM

[completar con tom y gtom](#)

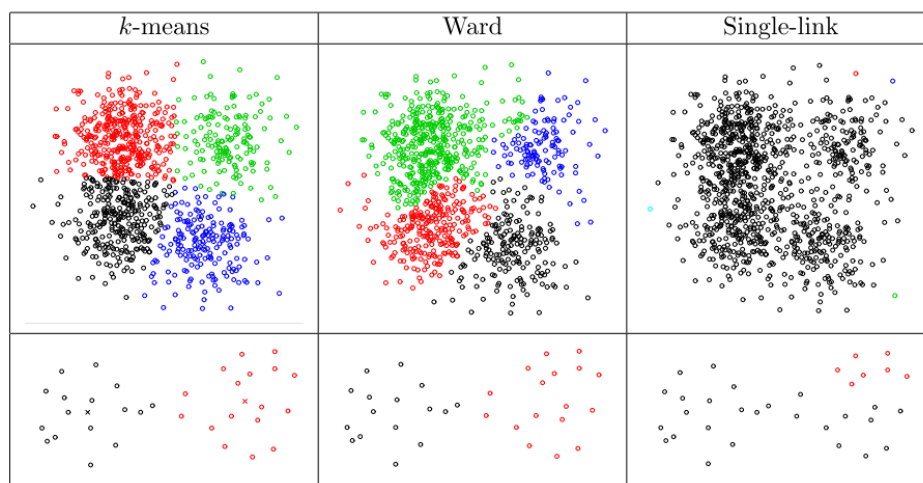


Figura 1.8: Algunos casos para los cuales los métodos de k -means o de Ward se comportan “mejor” que el método de enlace único (single-link). [citar correctamente](#)

caracterizacion de KTA local

Capítulo 2

conclusiones y perspectivas

Bibliografía

- [1] NATURE.COM. *Functional genomics*. Accedido: 2016-01-13.
URL <http://www.nature.com/subjects/functional-genomics>
- [2] WIKIPEDIA.ORG. *Functional genomics*. Accedido: 2016-01-13.
URL <https://en.wikipedia.org/wiki/Functional-genomics>
- [3] E. DOMANY. *Cluster Analysis of Gene Expression Data 1* **110** (2003) 1117.
- [4] B. ALBERTS. *Molecular Biology of The Cell*, volume 6 (2015).
- [5] B. BOSE. *In Vitro Differentiation of Pluripotent Stem Cells into Functional B Islets Under 2D and 3D Culture Conditions and In Vivo Preclinical Validation of 3D Islets*. *Methods in Molecular Biology* (2016) 257.
- [6] M. BABU. *An Introduction to Microarray Data Analysis*. *Computational Genomics: Theory and Application* (2004) 225.
URL <http://www.mrc-lmb.cam.ac.uk/genomes/madanm/microarray/chapter-final.pdf>
- [7] A. SCHULZE. *Navigating gene expression using microarrays: a technology review*. *Nature cell biology* **3** (2001) E190.
- [8] ARABIDOPSIS.ORG. *Microarray data from AtGenExpress*.
<https://www.arabidopsis.org/portals/expression/microarray/ATGenExpress.jsp>.
- [9] J. KILIAN *et al.* *The AtGenExpress global stress expression data set: Protocols, evaluation and model data analysis of UV-B light, drought and cold stress responses*. *Plant Journal* **50** (2007) 347.
- [10] ARABIDOPSIS-INTERACTOME-MAPPING-CONSORTIUM. *Evidence for Network Evolution in an Arabidopsis Interactome Map*. *Annual review of plant biology* **10** (2013) 161.
- [11] A. BRÜCKNER *et al.* *Yeast two-hybrid, a powerful tool for systems biology*. *International Journal of Molecular Sciences* **10** (2009) 2763.

- [12] M. E. CUSICK *et al.* *NIH Public Access*. Nature Methods **6** (2009) 39.
- [13] E. SEGAL *et al.* *Discovering molecular pathways from protein interaction and gene expression data*. Bioinformatics **19** (2003).
- [14] M. KANEHISA. *Yeast Biochemical Pathways. KEGG: Kyoto encyclopedia of genes and genomes*. Nucleic Acids Res **28** (2000) 27.
URL <http://pathway.yeastgenome.org/biocyc/>
- [15] J. PANDEY *et al.* *Functional coherence in domain interaction networks*. Bioinformatics **24** (2008) 28.
- [16] P. RESNIK. *Using Information Content to Evaluate Semantic Similarity in a Taxonomy*. roceedings of the 14th international joint conference on Artificial intelligence - Volume 1 - IJCAI'95 **1** (1995) 6.
URL <http://arxiv.org/abs/cmp-lg/9511007>
- [17] C. PESQUITA *et al.* *Semantic similarity in biomedical ontologies*. PLoS Computational Biology **5** (2009).
- [18] A. BERENSTEIN. *Análisis de redes complejas en sistemas biomoleculares* (2014).
- [19] ASHBURNER. *Gene ontology: tool for the unification of biology*. Nat Genet **25** (2000).
- [20] G. GAN *et al.* *Data Clustering: Theory, Algorithms, and Applications*, volume 20 (2007).
- [21] M. HALKIDI *et al.* *On clustering validation techniques*. Journal of Intelligent Information Systems **17** (2001) 107.
- [22] E. DOMANY. *Superparamagnetic clustering of data—the definitive solution of an ill-posed problem*. Physica A: Statistical Mechanics and its Applications **263** (1999) 158.
URL <http://www.sciencedirect.com/science/article/pii/S0378437198004944>
- [23] S. CHEN *et al.* *On the similarity metric and the distance metric*. Theoretical Computer Science **410** (2009) 2365.
URL <http://dx.doi.org/10.1016/j.tcs.2009.02.023>
- [24] L. W. KHENG. *Image Registration* (2010).
- [25] P. D'HAESELEER. *How does gene expression clustering work?* Nat Biotech **24** (2005).

- [26] C. HENNIG. *How to find an appropriate clustering for mixed-type variables with application to socio-economic stratification*. Journal of the Royal Statistical Society. Series C: Applied Statistics **62** (2013) 309.
- [27] M. EISEN. *Cluster analysis and display of genome-wide expression patterns*. Proceedings of the National Academy of Sciences of the United States of America **95** (1998) 14863.
- [28] D. LIN. *An Information-Theoretic Definition of Similarity*. In: Proc. of the 15th International Conference on Machine Learning (1998) 296.
- [29] J. JIANG. *Semantic Similarity Based on Corpus Statistics and Lexical Taxonomy*. Proceedings of International Conference Research on Computational Linguistics (1997) 19.
- [30] H. K. LEE *et al.* *Coexpression Analysis of Human Genes Across Many Microarray Data Sets* (2004) 1085.
- [31] J. SEVILLA. *Correlation between gene expression and go semantic similarity*. In: IEEE/ACM Transactions on Computational Biology and Bioinformatics (2005).
- [32] P. LORD. *Investigating semantic similarity measures across the Gene Ontology: the relationship between sequence and annotation*. Bioinformatics (2003).
- [33] J. KOGAN. *Introduction to Clustering Large and High-Dimensional Data* (2006).
- [34] J. HARTIGAN. *A K-Means Clustering Algorithm*. Journal of the Royal Statistical Society **28** (1979) 100.
- [35] H. S. PARK. *A simple and fast algorithm for K-medoids clustering*. Expert Systems with Applications **36** (2009) 3336.
- [36] L. IBRAHIM. *Using Modified Partitioning Around Medoids Clustering Technique in Mobile Network Planning* **9** (2012) 299.
- [37] J. STEPHEN. *Hierarchical clustering schemes*. Psychometrika (1967).
- [38] C. SHALIZI. *Distances between Clustering , Hierarchical Clustering*. Data Mining (2009) 36.
- [39] P. LANGFELDER *et al.* *Defining clusters from a hierarchical cluster tree: The Dynamic Tree Cut package for R*. Bioinformatics **24** (2008) 719.
- [40] P. LANGFELDER *et al.* *Dynamic Tree Cut : in-depth description , tests and applications* (2007) 1.

- [41] M. ROSVALL. *Maps of random walks on complex networks reveal community structure*. Proceedings of the National Academy of Sciences of the United States of America **105** (2008) 1118.