



Redes y Sistemas Complejos
Cuarto Curso del Grado en Ingeniería Informática

Tema 6: Modularidad, Particionamiento y Comunidades

Oscar Cordon García

Dpto. Ciencias de la Computación e Inteligencia Artificial
ocordon@decsai.ugr.es

ESTRUCTURA DE COMUNIDADES

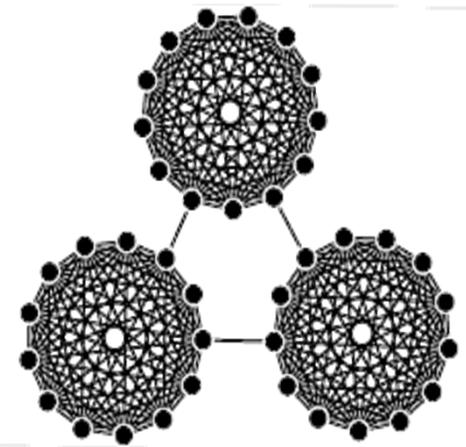
Las redes complejas tienden a mostrar una **estructura de comunidades**

Esta propiedad suele darse como consecuencia de la heterogeneidad global y local de la distribución de los enlaces en un grafo

A menudo encontramos una **alta concentración de enlaces en ciertas regiones del grafo**, denominadas **comunidades**, y una **baja concentración de enlaces entre esas regiones**

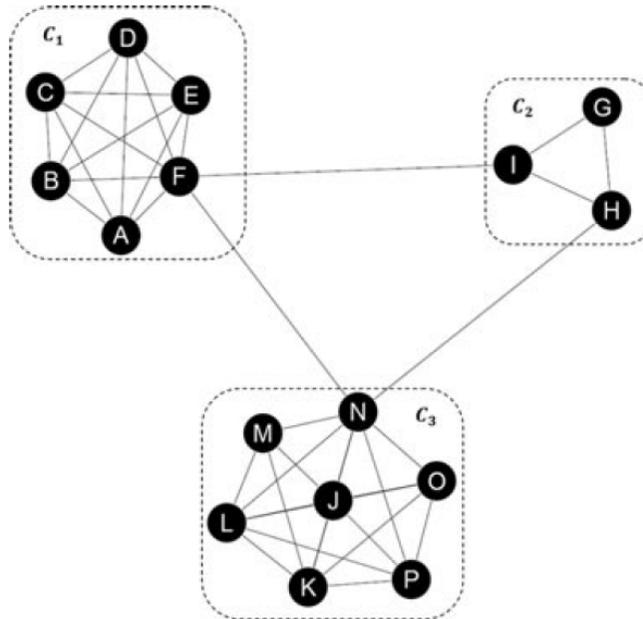
Las comunidades, también conocidas como módulos o clusters, se definen de forma sencilla como grupos de nodos similares

A partir del concepto de densidad de la red, **las comunidades pueden definirse como grupos de nodos densamente conectados que presentan conexiones dispersas entre sí**



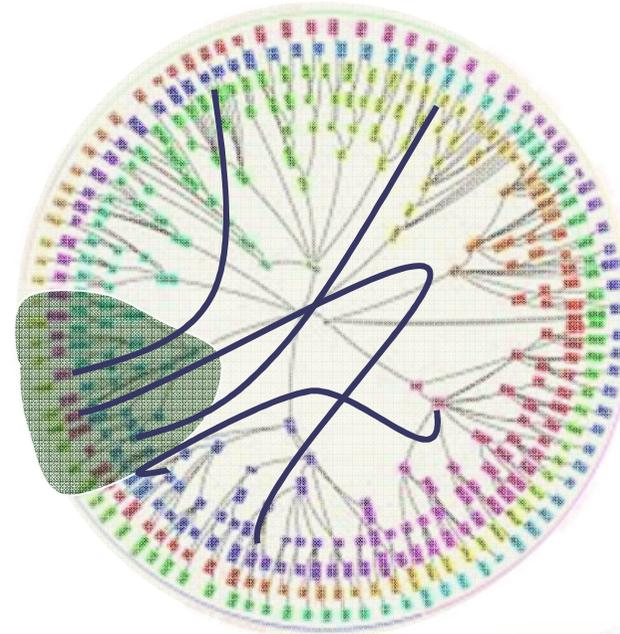
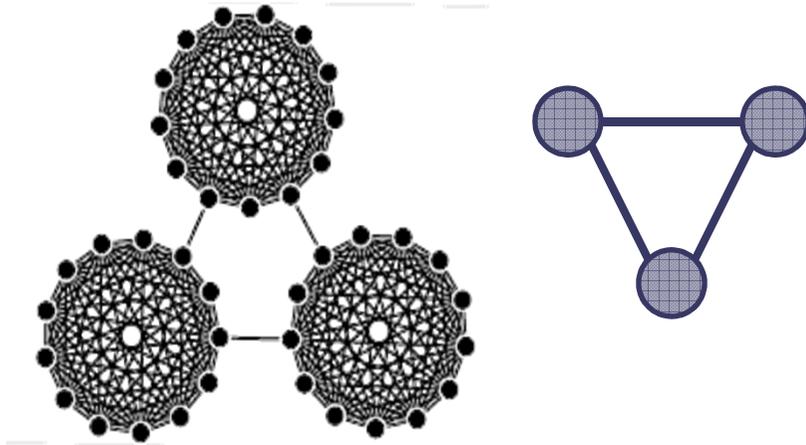
INTRODUCCIÓN

Caso ideal



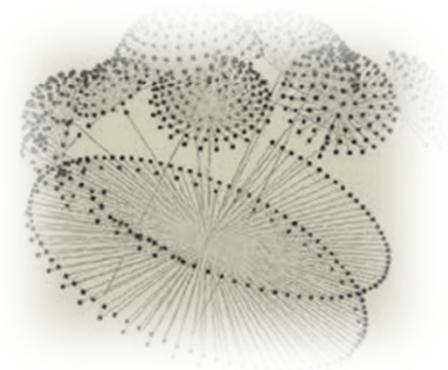
- En esta red, hay **tres comunidades**: C_1 , C_2 y C_3
- Cada comunidad está formada por un grafo completo (un **clique**) de tamaño variable ($C_1 = K_6$, $C_2 = K_3$ y $C_3 = K_7$)
- La densidad de enlaces entre las comunidades es muy baja. Los pocos enlaces que existen son **puentes**

El clustering implica modularidad

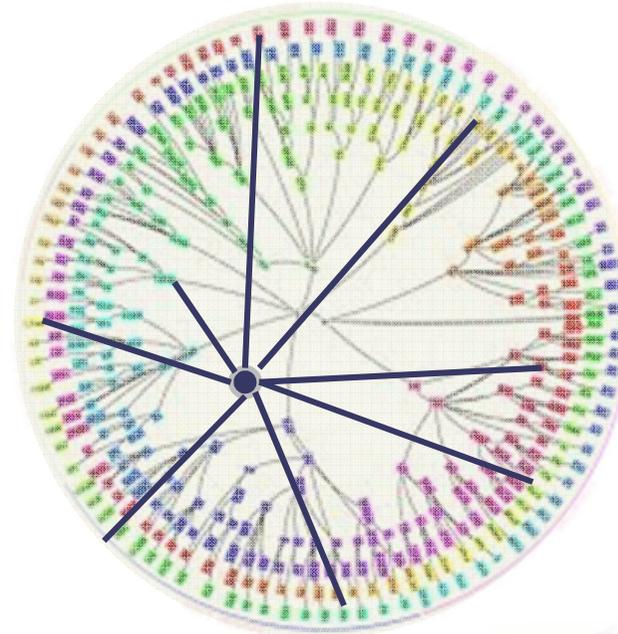
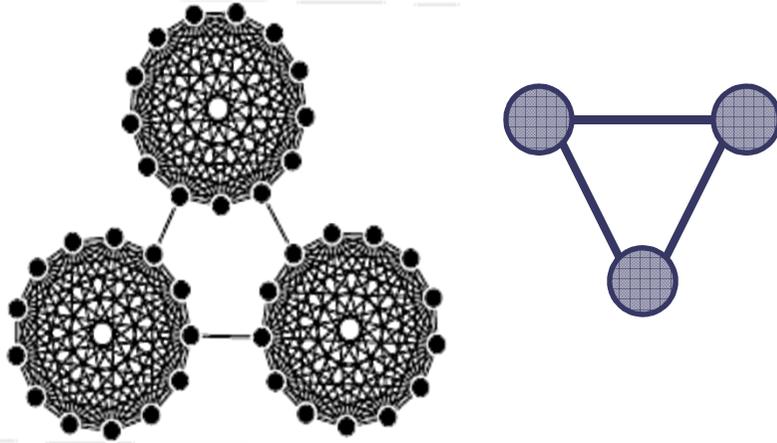


La funcionalidad requiere modularidad

La propiedad de mundos pequeños tiende a eliminar la modularidad



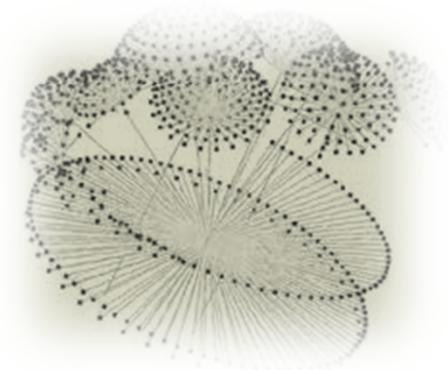
El clustering implica modularidad



La funcionalidad requiere modularidad

La propiedad de mundos pequeños tiende a eliminar la modularidad

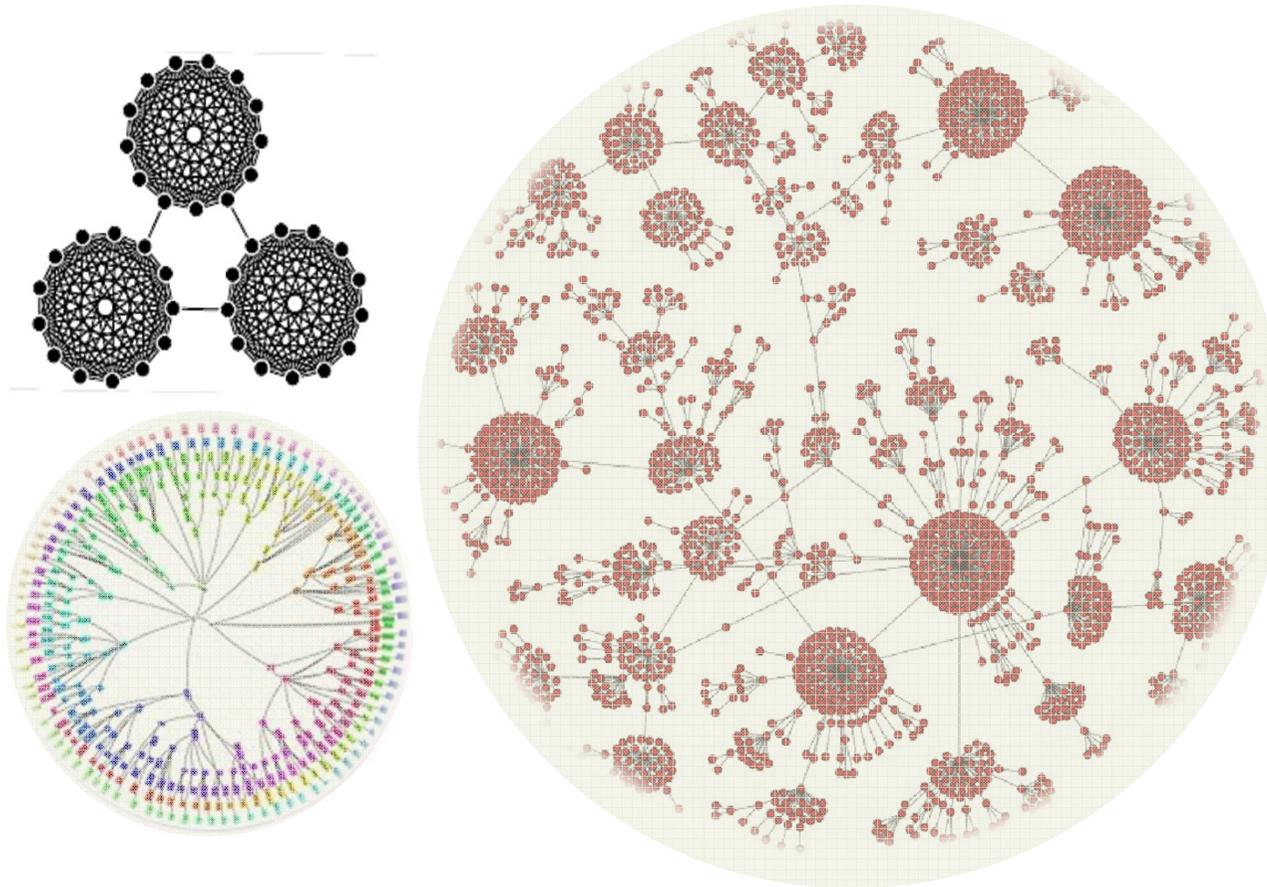
Los hubs tienden a eliminar la modularidad



ESTRUCTURA MODULAR DE LAS REDES

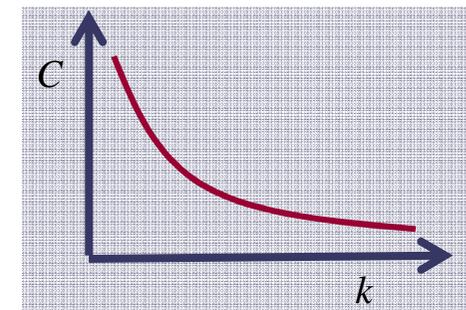
¿Cuándo es modular una red?

El clustering sólo debe estar en la periferia



Los nodos de grado bajo suelen pertenecer a un único módulo

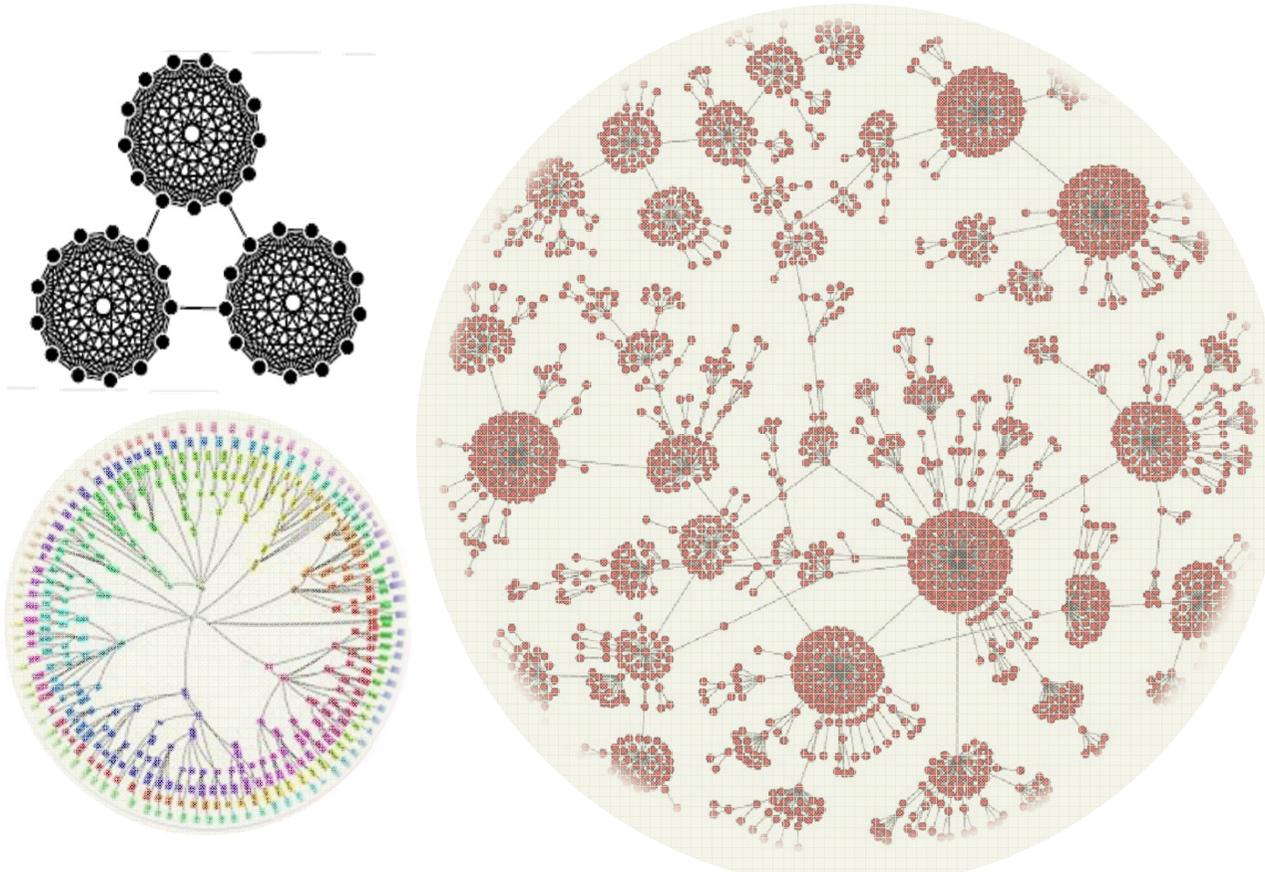
Los hubs hacen de puente entre distintos módulos



ESTRUCTURA MODULAR DE LAS REDES

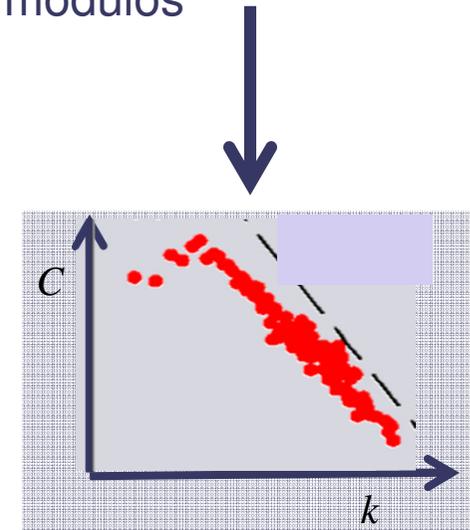
¿Cuándo es modular una red?

El clustering sólo debe estar en la periferia



Los nodos de grado bajo suelen pertenecer a un único módulo

Los hubs hacen de puente entre distintos módulos



Pero... ¿cómo descubrimos los módulos?

La modularidad Q es una función de calidad que mide la calidad de una partición concreta de una red en comunidades:

Se define como la diferencia entre el número de enlaces existentes en los grupos y el número de enlaces esperado en una red aleatoria equivalente

$$Q = \frac{1}{2L} \sum_{ij} \left[A_{ij} - \frac{k_i k_j}{2L} \right] \delta(c_i, c_j)$$

número de enlaces de la red \rightarrow $2L$
 matriz de adyacencia \rightarrow A_{ij}
 la probabilidad de un enlace entre dos nodos es proporcional a sus grados \rightarrow $\frac{k_i k_j}{2L}$
 vale 1 si los nodos son de la misma comunidad \rightarrow $\delta(c_i, c_j)$

La idea básica es que la red muestra una estructura modular coherente si el número de enlaces entre comunidades es menor que el esperado en una red aleatoria

$Q \in [-1,1]$. Cuanto mayor es su valor, mejor es la partición, es decir, las comunidades encontradas están densamente conectadas internamente (hay más enlaces de los que cabría esperar aleatoriamente) y dispersamente conectadas entre sí

En una red aleatoria, $Q=0$. En la práctica, una modularidad de 0.3 es un buen valor

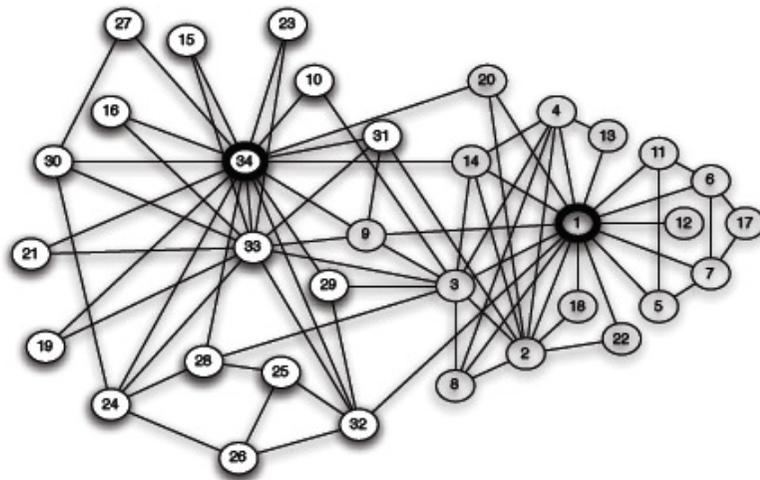
Se usa tanto para comparar la calidad de distintas particiones como diseñar métodos de descubrimiento de comunidades que traten de maximizar su valor

¿PARA QUÉ DETECTAR
COMUNIDADES?

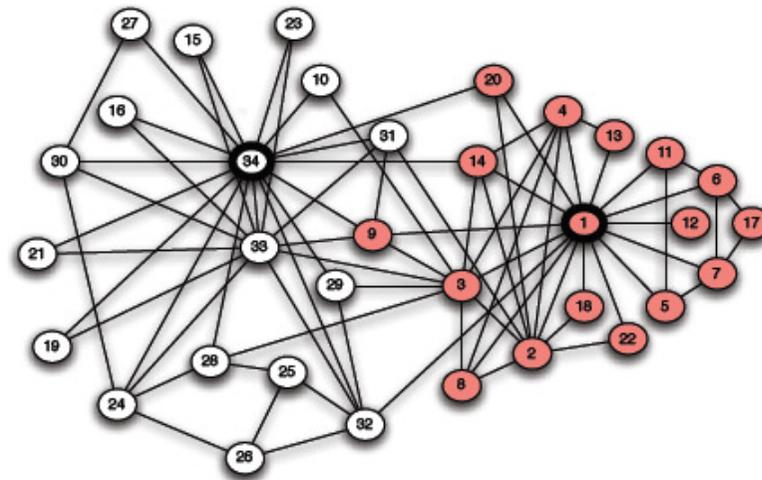
COMUNIDADES EN EL MUNDO REAL (1)

En la vida real existen muchos ejemplos de grupos compactos en redes complejas:

- Sociedades: las personas tienen una tendencia natural a formar grupos (familias, círculos de amigos, grupos profesionales o religiosos, ciudades, naciones, etc.)
- Empresariales/Económicas: compañías, clientes, etc.
- Biología: p.ej. redes metabólicas. En redes de interacción de proteínas encontramos grupos de proteínas con funciones similares dentro de la célula
- Internet: comunidades virtuales (Facebook, Twitter, etc.), grupos de páginas web relacionadas, etc. (útiles para sistemas de recomendaciones)
- y muchos ámbitos más...



(a) *Karate club network*

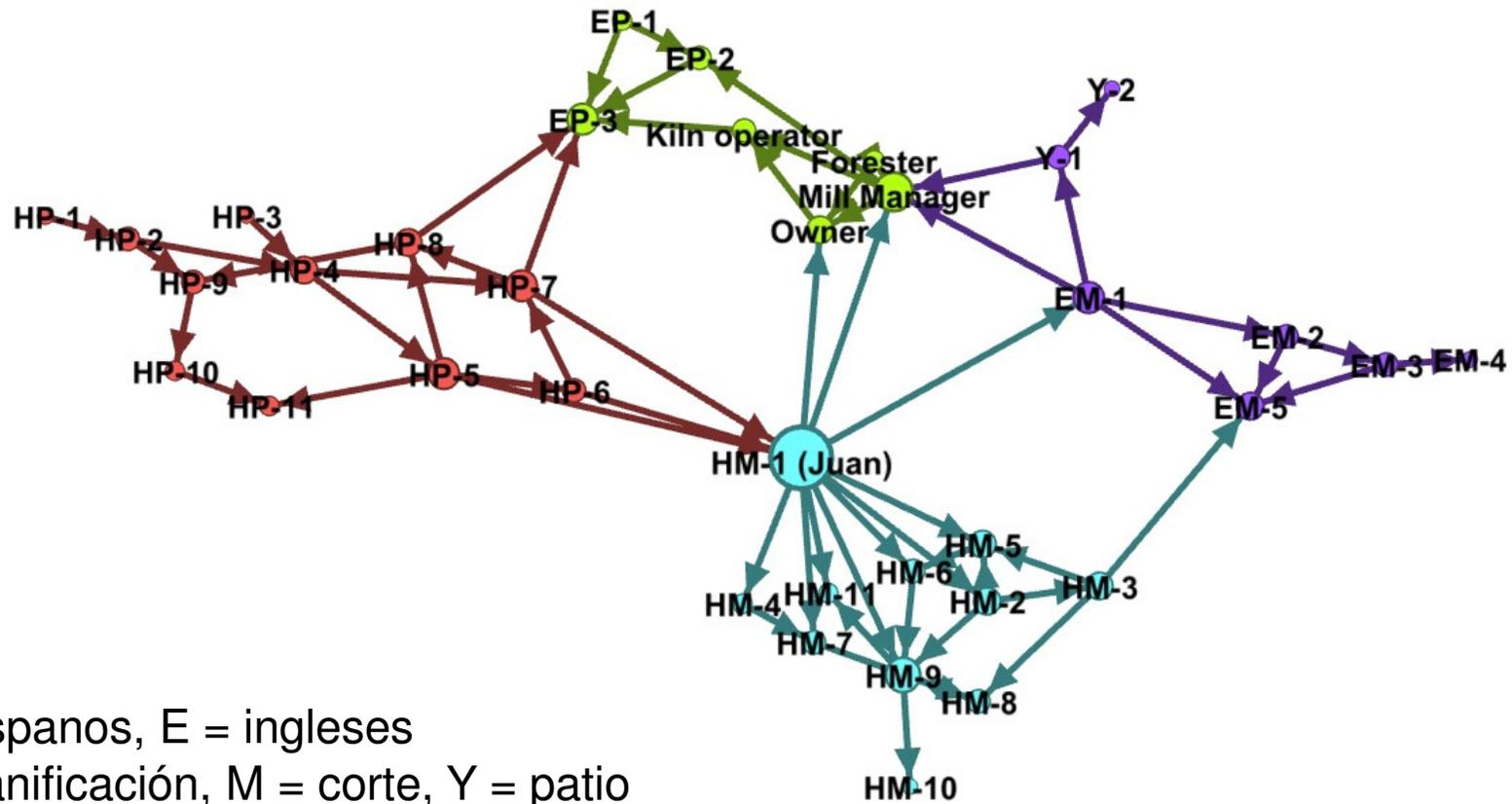


(b) *After a split into two clubs*

El club de karate de Zachary

Zachary. An information flow model for conflict and fission in small groups. *J. Anthropol. Res.* 33: 452-473 (1977)

FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (1)



H = hispanos, E = ingleses
P = planificación, M = corte, Y = patio

Red del aserradero (*sawmill*). Fuente: Exploratory Social Network Analysis with Pajek

FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (2)

La dirección del aserradero estaba teniendo dificultades para persuadir a los trabajadores para que adoptaran un nuevo plan que redundaría en un beneficio para todos los empleados

En concreto, los trabajadores hispanos eran los más reticentes a aceptar

La dirección contrató a un sociólogo que diseñó una red que reflejaba con qué compañeros hablaban habitualmente

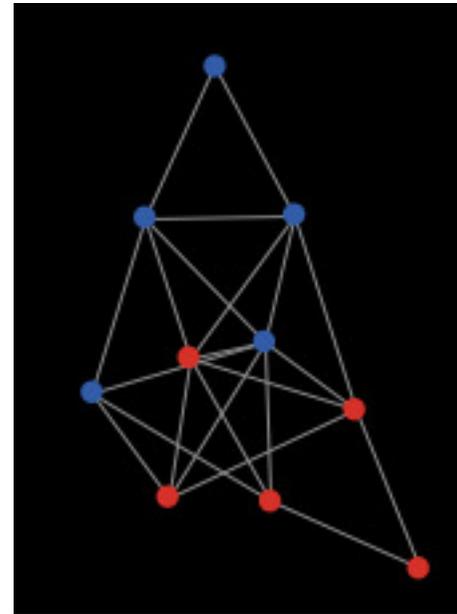
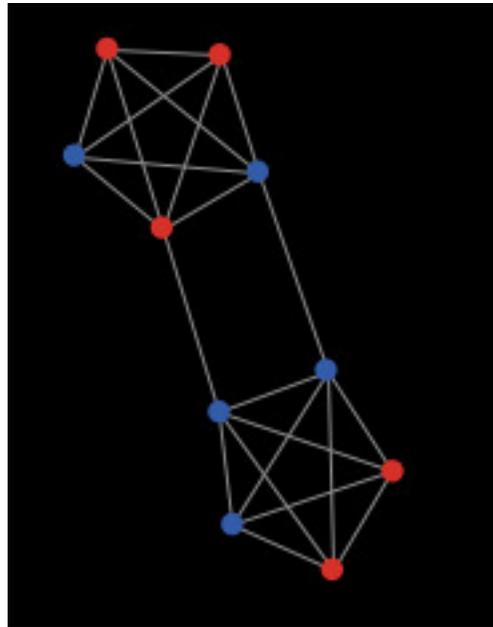
El sociólogo recomendó a la dirección que hablaran con Juan y que le pidieran que hablara con el resto de trabajadores hispanos

Fue un éxito, rápidamente todos estaban de acuerdo con el nuevo plan

¿Por qué?

FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (3)

- <http://www.ladamic.com/netlearn/NetLogo502/OpinionFormationModelToy.html>
- Cada nodo adopta la opinión mayoritaria de sus vecinos (una opinión aleatoria, en caso de empates)



FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (4)

PREGUNTA:

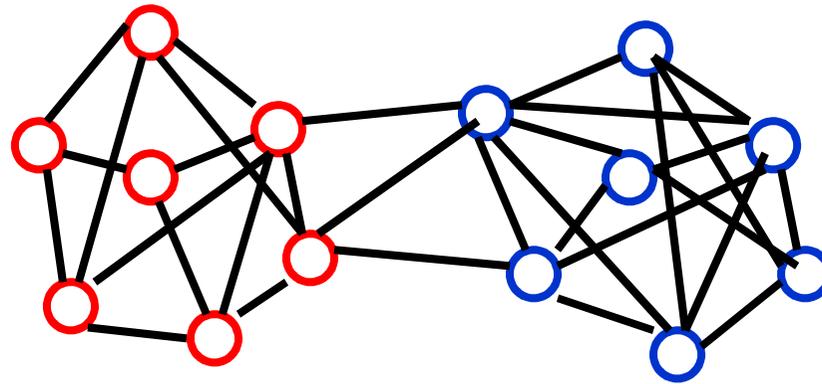
Ejecutar alternativamente la configuración de dos comunidades y la de Erdos-Renyi. ¿Cuál puede mantener opiniones divergentes cuando se itera la actualización de opiniones?

- a) Sólo la de Erdos-Renyi
- b) Sólo la de dos comunidades
- c) Ambas

<http://www.ladamic.com/netlearn/NetLogo502/OpinionFormationModelToy.html>

FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (5)

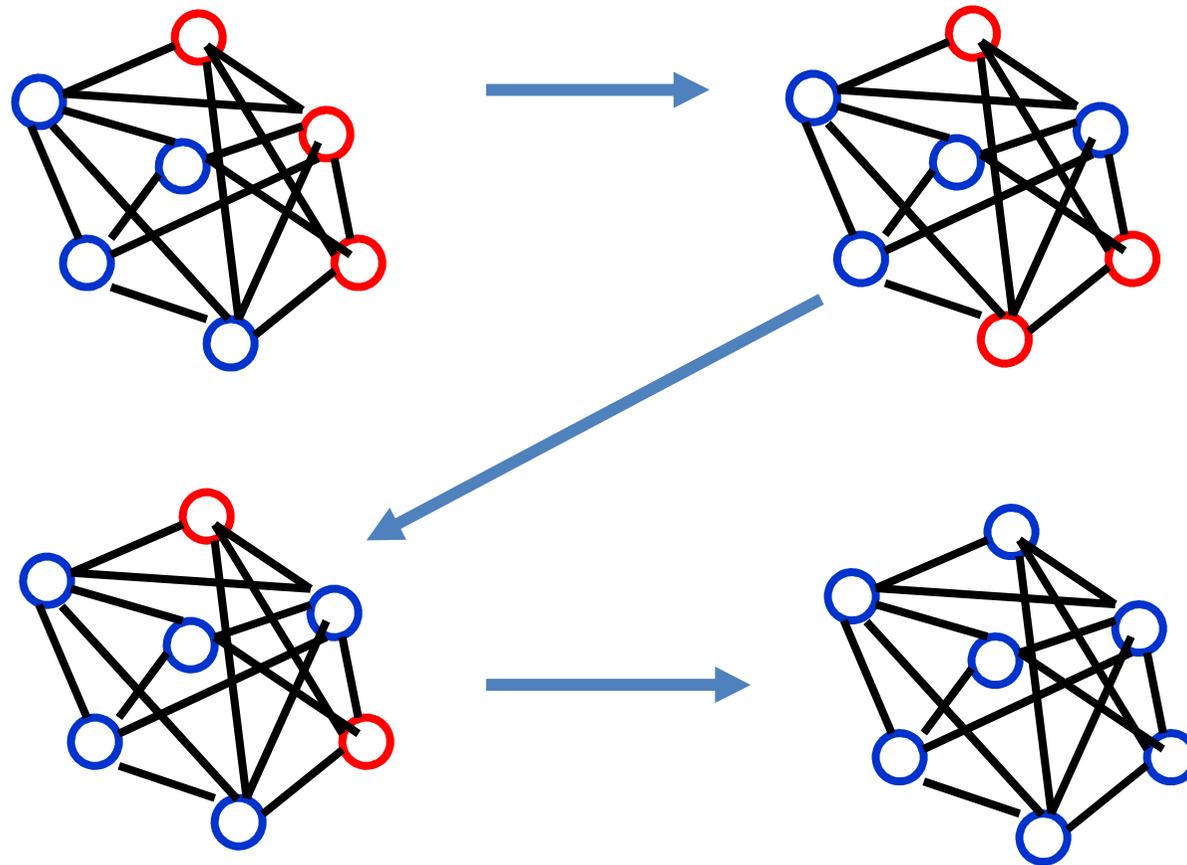
FORMACIÓN DE OPINIONES Y UNIFORMIDAD:



Si cada nodo adopta la opinión de la mayoría de sus vecinos, es posible formar opiniones distintas en subgrupos cohesivos distintos

FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (6)

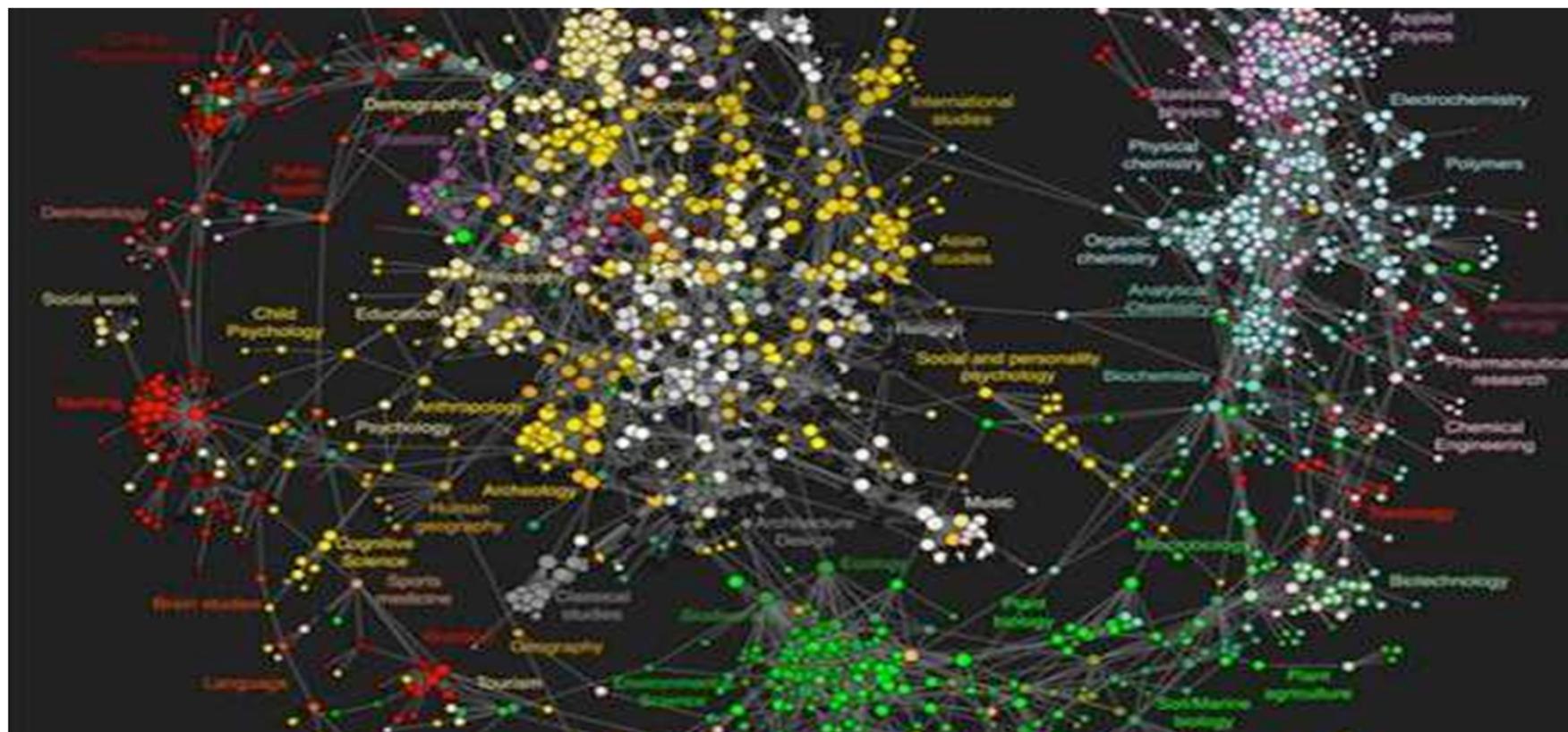
Hay más uniformidad dentro de un grupo cohesivo:



FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (7)

Rosvall y Bergstrom. Maps of random walks on complex networks reveal community structure. Proc. Natl. Acad. Sci. USA 105: 1118-1123 (2008)

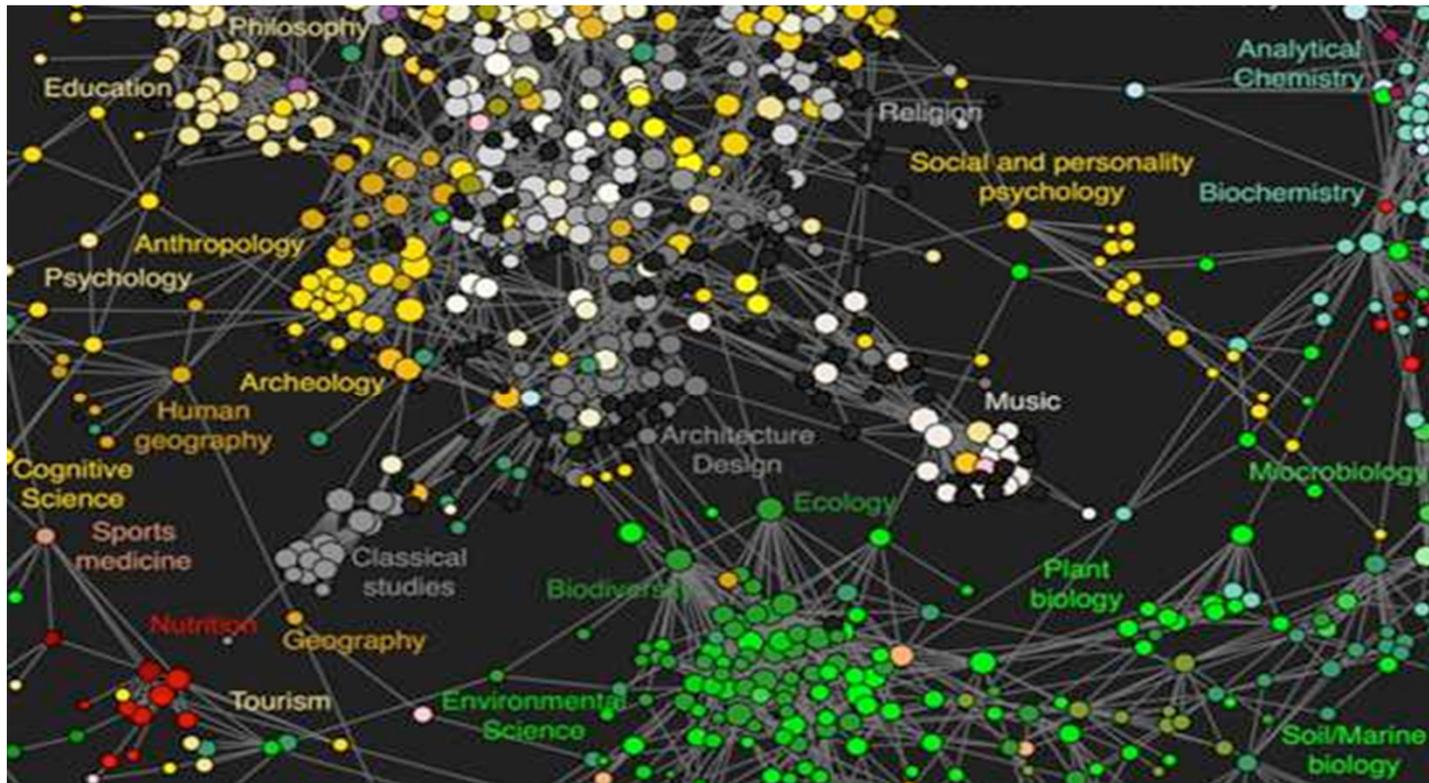
Mapa de la ciencia de citas de 6000 revistas



FORMACIÓN DE OPINIONES Y ESTRUCTURA DE COMUNIDADES (8)

Rosvall y Bergstrom. Maps of random walks on complex networks reveal community structure. Proc. Natl. Acad. Sci. USA 105: 1118-1123 (2008)

Clustering para comprimir la información de *random walks* que siguen el flujo de las citas de un campo científico a otro, de modo que las áreas emergen de forma natural



MÉTODOS DE DETECCIÓN DE COMUNIDADES

Fortunato. Community detection in graphs. Phys Rep 486: 75–174 (2010)

CRITERIOS ESTRUCTURALES QUE IDENTIFICAN UNA COMUNIDAD EN UNA RED

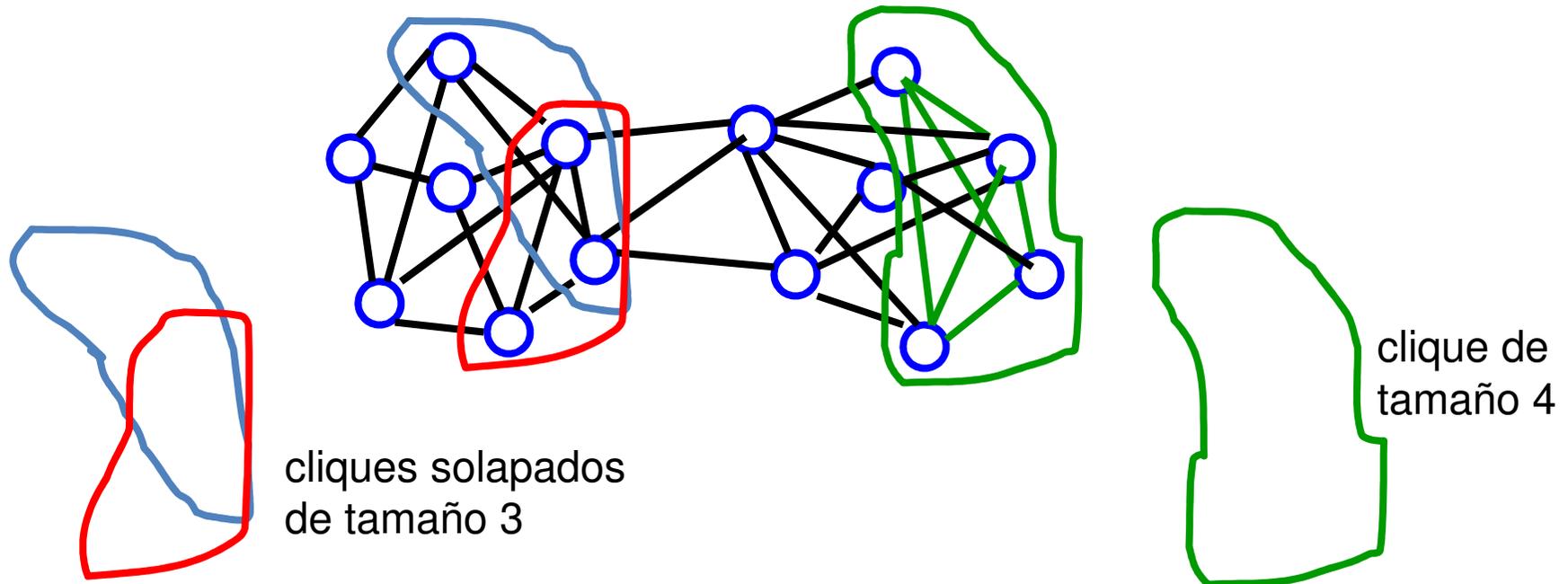
- **Mutualidad completa** (*cliques*):
 - El grupo es un subgrafo completo (todo el mundo se conoce en el grupo)
- **Frecuencia de enlaces entre los miembros** (*k-cores*):
 - Todos los miembros del grupo tiene enlaces al menos a otros k miembros
- **Alcanzabilidad/cercanía entre los miembros** (*n-cliques*):
 - Los individuos del grupo están separados por un máximo de q saltos
- **Comparación de la cohesión interna y externa del grupo** (*p-cliques*):
 - Frecuencia relativa de enlaces entre los miembros del grupo en comparación con la de los no miembros

Wasserman y Faust. *Social Network Analysis*. Cambridge University Press; 1994

IDENTIFICACIÓN DE CLIQUES

Concepto

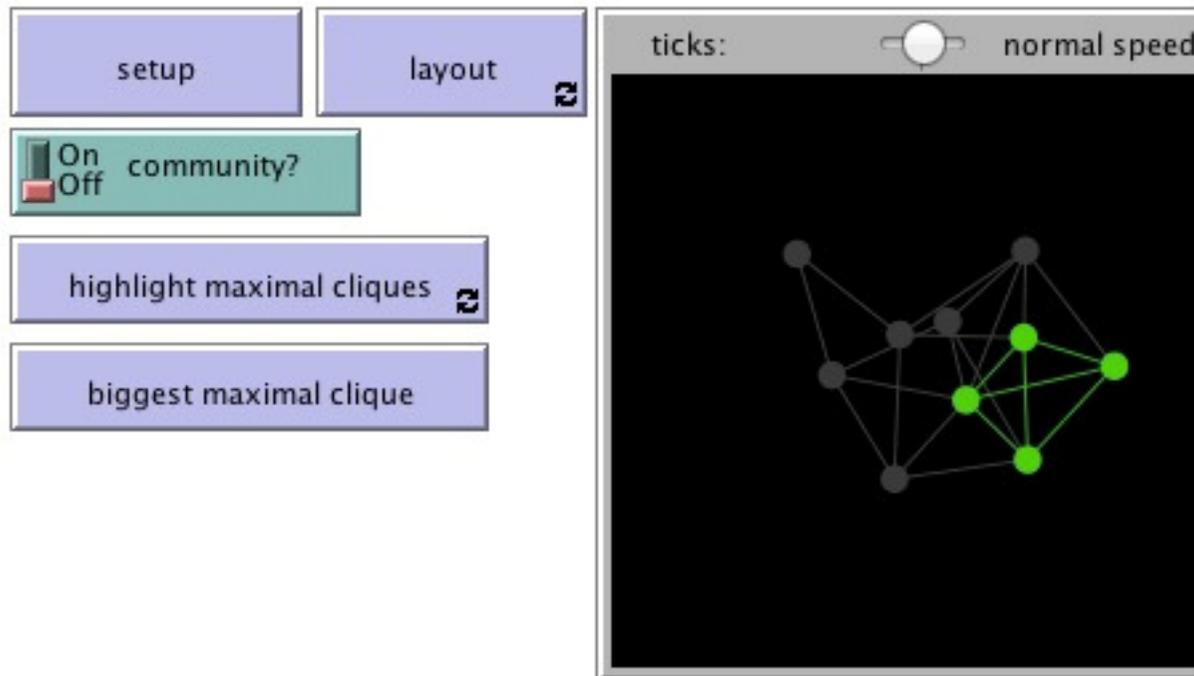
- Todos los miembros del grupo tienen enlaces al resto
- Los triángulos son los cliques más básicos, los de mayor dimensión son menos frecuentes
- Los cliques pueden estar solapados



IDENTIFICACIÓN DE CLIQUES

Modelo Netlogo (1)

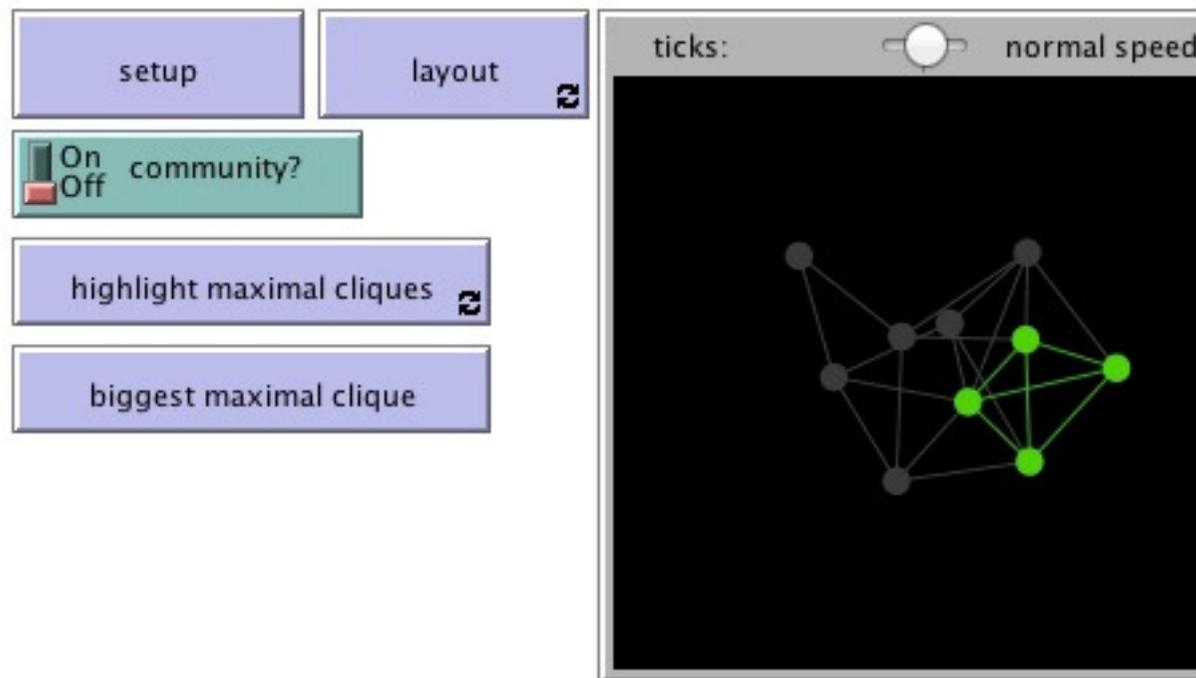
- <http://www.ladamic.com/netlearn/nw/Cliques.html>
- Comparar la configuración de red aleatoria ER con la de estructura de comunidades (son las mismas que en el modelo de formación de opiniones)



IDENTIFICACIÓN DE CLIQUES

Modelo Netlogo (2)

- **PREGUNTA:** ¿Cuál de los dos tiene un clique máximo de mayor tamaño?
- Los cliques revelan la estructura de comunidades...

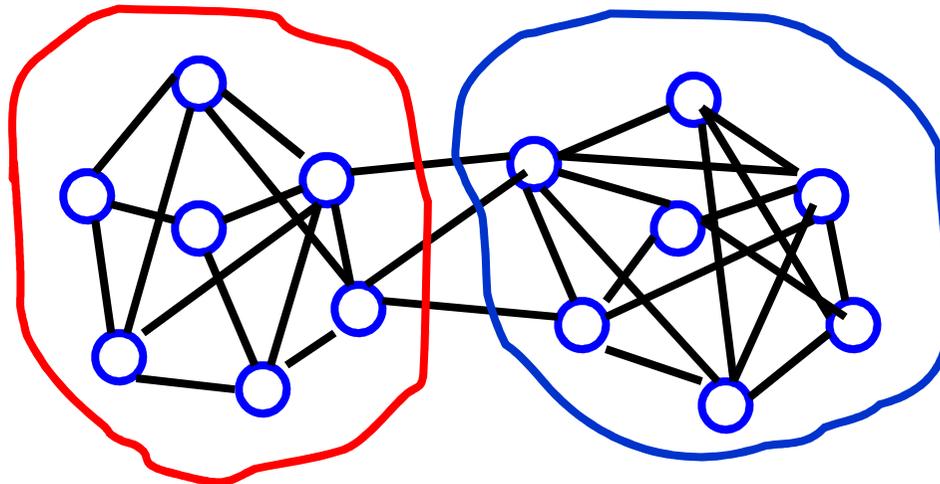


- Es un problema NP-completo
- No son robustos
 - Un solo enlace faltante puede descalificar un clique, haciendo que el grupo no sea considerado una comunidad
- No son interesantes
 - Todo el mundo está conectado entre sí
 - No hay estructura centro-periferia, no hay jerarquía entre los enlaces
 - Las medidas de centralidad no dan información
- El solapamiento de los cliques puede ser más relevante que su propia existencia...

IDENTIFICACIÓN DE K-CORES

Frecuencia de enlaces entre miembros (1)

- Cada nodo de un grupo está conectado con otros k nodos de dicho grupo:



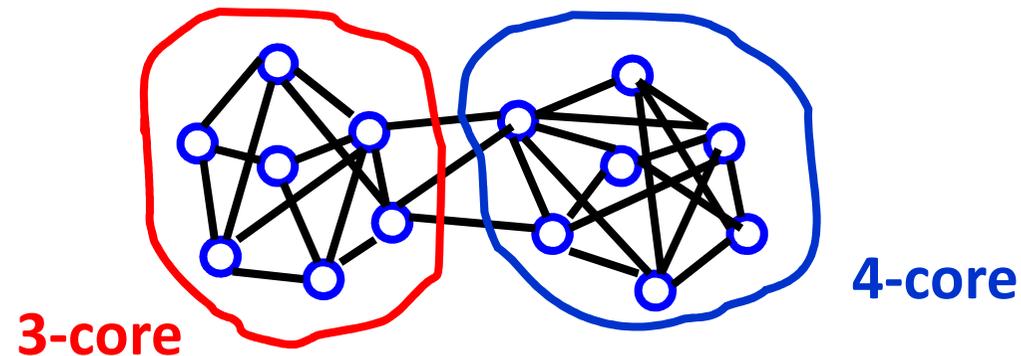
Pregunta:

- ¿Cuál es el valor de k para el k -core marcado en rojo?
- ¿y para el azul?

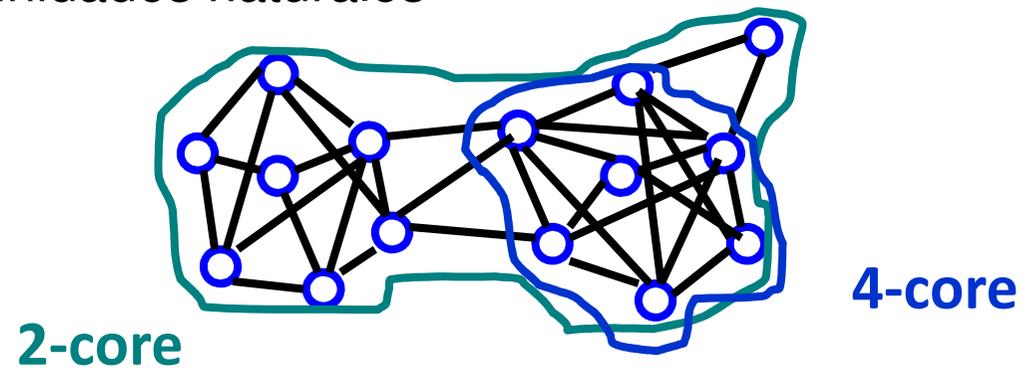
IDENTIFICACIÓN DE K-CORES

Frecuencia de enlaces entre miembros (2)

- Cada nodo de un grupo está conectado con otros k nodos de dicho grupo:



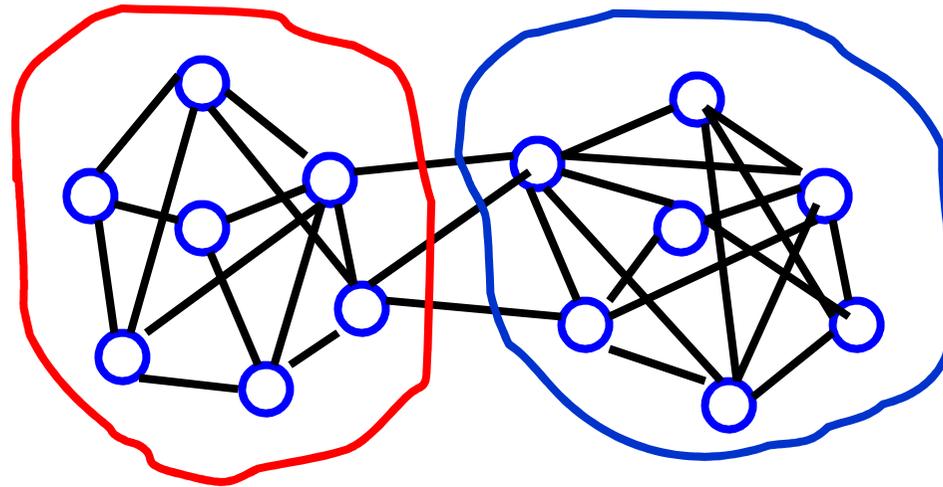
- Aun así, es una estructura demasiado restrictiva como requisito para identificar comunidades naturales



IDENTIFICACIÓN DE N-CLIQUES

Alcanzabilidad y diámetro

- La máxima distancia entre cualesquiera dos nodos del grupo es n (el 1-clique equivale al clique):



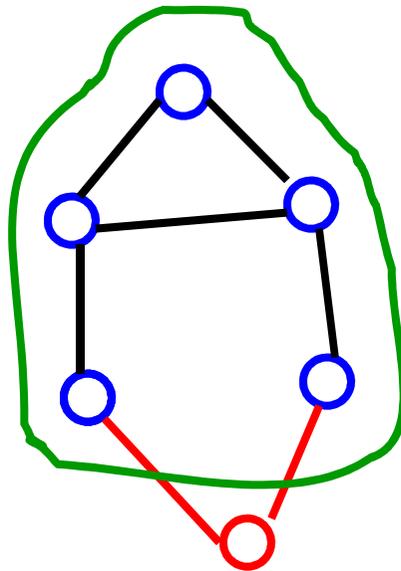
2-cliques

- Justificación teórica: Flujo de información a través de intermediarios

IDENTIFICACIÓN DE N-CLIQUE

Consideraciones

- Problemas:
 - El diámetro puede ser mayor que n
 - El n -clique puede estar desconectado (los caminos pueden pasar por nodos que no estén en el grupo)



2-clique
diámetro = 3

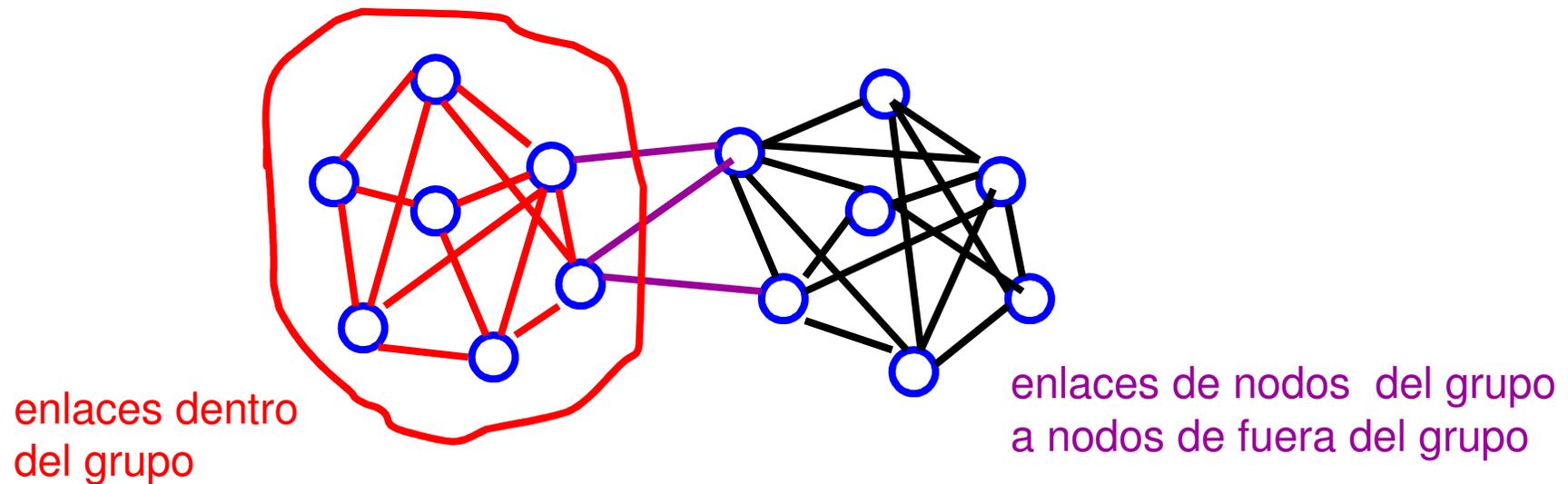
camino externo al 2-clique

- **Solución:** n -clubs: subgrafos máximos de diámetro 2

IDENTIFICACIÓN DE P-CLIQUES

Cohesión de enlaces en el grupo

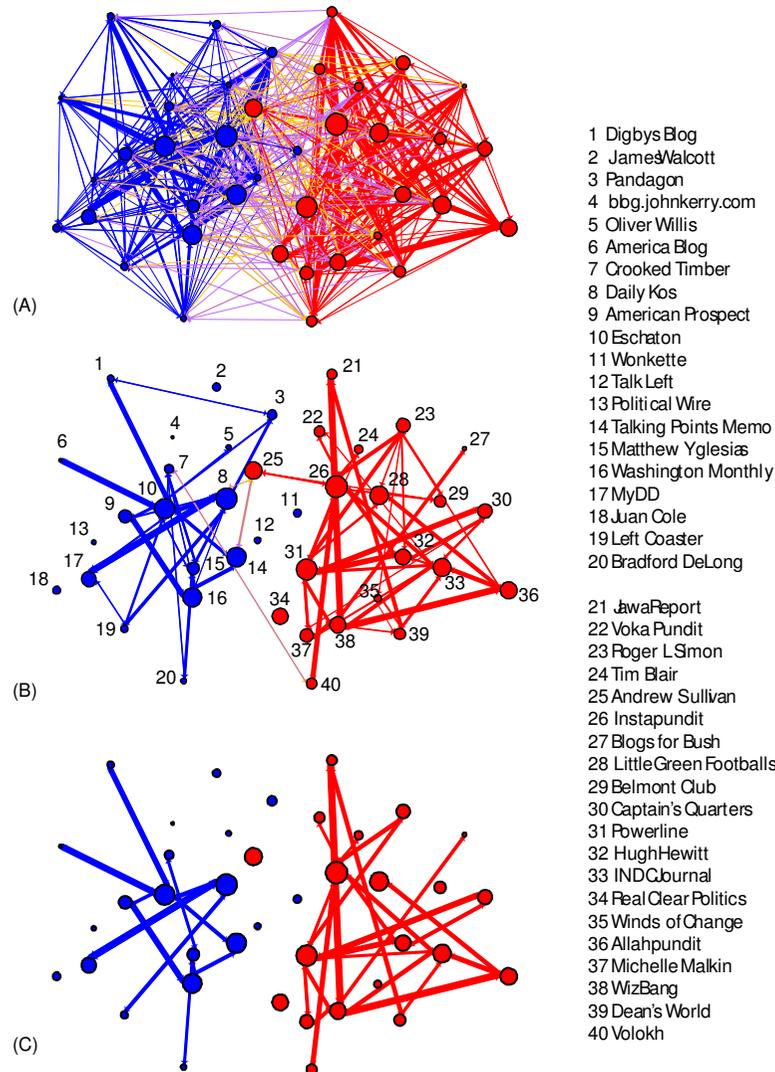
- Particionamiento de la red en clusters en los que los nodos tienen como mínimo una proporción $p \in [0,1]$ de vecinos dentro del grupo:



- Se aplica obteniendo **componentes fuertemente conexas**
- Es un proceso costoso. Es conveniente podar enlaces antes de aplicarlo:
 - Determinar reciprocidades
 - Poda de enlaces basada en un umbral sobre los pesos

COHESIÓN EN REDES DIRIGIDAS Y PONDERADAS

Ejemplo



Ejemplo: Blogs políticos
(29 Ago-15 Nov 2004)

- A. Todas las citas entre blogs en los dos meses anteriores a la elección de 2004
- B. Citas entre blogs con al menos 5 citas en ambas direcciones
- C. Poda más restrictiva incluyendo enlaces con al menos 25 citas combinadas

Sólo un 15% de las citas construyen comunidades

Adamic y Glance. The Political Blogosphere and the 2004 U.S. Election: Divided They Blog. Proc. LinkKDD2005

ENFOQUES PARA DETECCIÓN DE COMUNIDADES

Según Newman y Girvan (2004), existen dos líneas de investigación principales para el descubrimiento de comunidades en redes complejas:

1. **Particionamiento de grafos**: tiene su origen en Informática, en el campo de la computación distribuida. Busca la mejor forma de asignar tareas a procesadores para minimizar las comunicaciones entre ellos
2. **Modelado de bloques** (también llamado **clustering jerárquico** o **detección de la estructura de comunidades**): se origina en Sociología. Está motivado por el descubrimiento de grupos en una sociedad para facilitar el análisis de fenómenos sociales

En cualquier caso, el procedimiento implica dividir el grafo original en un conjunto de subgrafos disjuntos mediante la optimización de una función objetivo (p.ej. la modularidad)

Newman y Girvan. Finding and evaluating community structure in networks. Phys Rev E 69:026113 (2004)

CARACTERÍSTICAS GENERALES

- El propósito de los dos enfoques es descubrir grupos de nodos relacionados en la red y, si es posible, la estructura jerárquica correspondiente, a partir de la información proporcionada por la topología de la red
- Una de las heurísticas más extendidas es eliminar iterativamente los puentes entre grupos de nodos (Girvan y Newman, 2002)
- Estos métodos devuelven particiones disjuntas del conjunto de nodos, cada nodo pertenece a una única comunidad (no permiten el solapamiento de comunidades)
- Existen algunos que consideran dicho solapamiento, como el *clique percolation method* de Palla y otros (2005). Son especialmente útiles en Ciencias Sociales

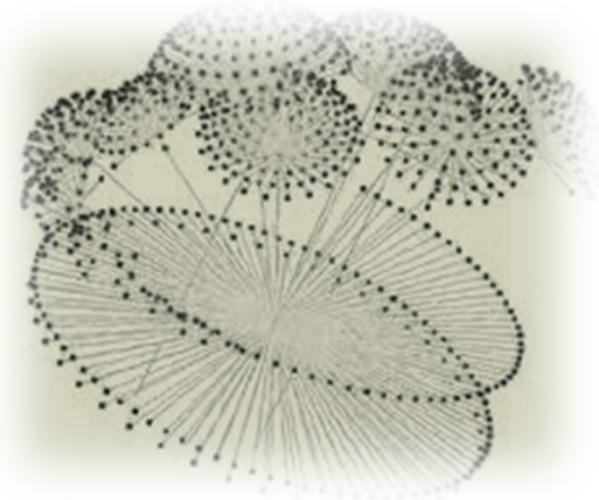
Girvan y Newman. Community structure in social and biological networks. Proc Natl Acad Sci USA 99:7821–7826 (2002)

Palla y otros. Uncovering the overlapping community structure of complex networks in nature and society. Nature 435:814–818 (2005)

PARTICIONAMIENTO DE GRAFOS/REDES COMPLEJAS (1)

División óptima de la red en un **número predeterminado** de particiones

$$\left\{ \begin{array}{l} x_1 = \dots \\ x_2 = \dots \\ \vdots \\ x_n = \dots \end{array} \right.$$

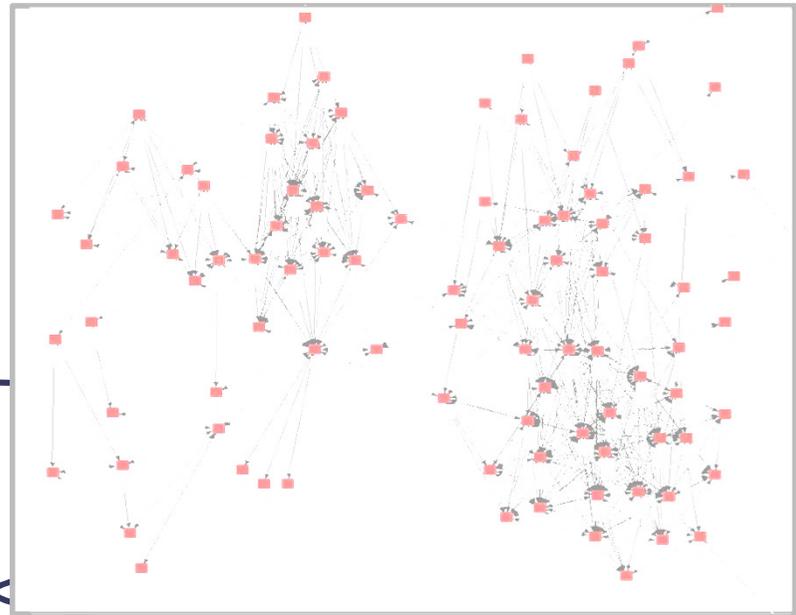
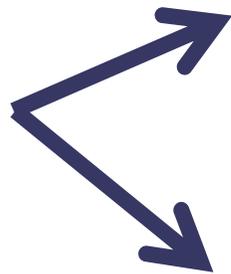


División de una tarea en sub-tareas

PARTICIONAMIENTO DE GRAFOS/REDES COMPLEJAS (2)

División óptima de la red en un **número predeterminado** de particiones

$$\left\{ \begin{array}{l} x_1 = \dots \\ x_2 = \dots \\ \vdots \\ x_n = \dots \end{array} \right.$$

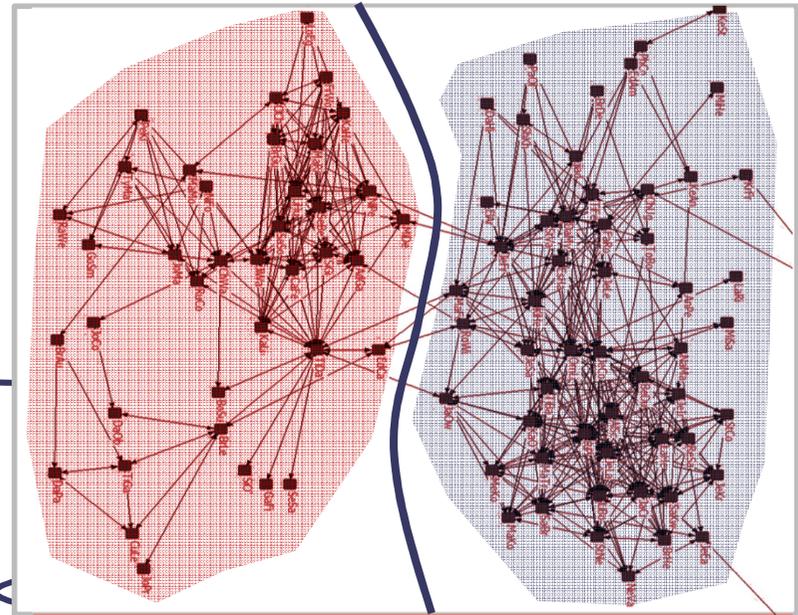
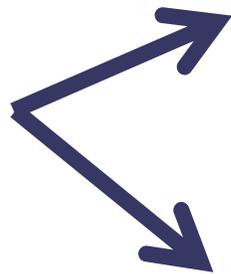


División de una tarea en sub-tareas, minimizando la comunicación entre tareas

PARTICIONAMIENTO DE GRAFOS/REDES COMPLEJAS (3)

División óptima de la red en un **número predeterminado** de particiones

$$\left\{ \begin{array}{l} x_1 = \dots \\ x_2 = \dots \\ \vdots \\ x_n = \dots \end{array} \right.$$



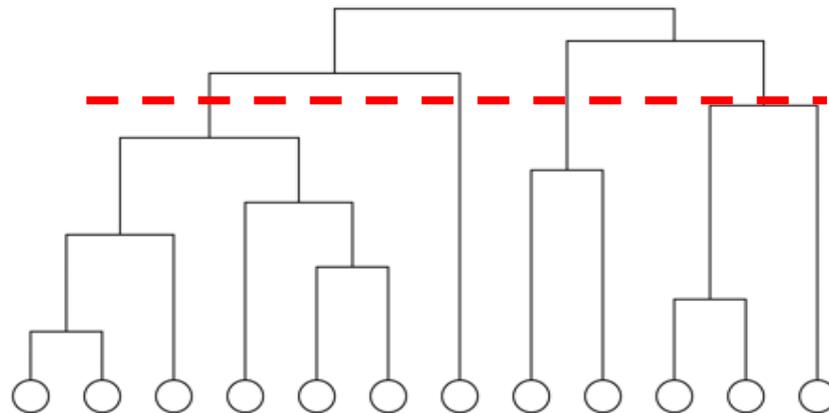
En redes, el número y el tamaño de las comunidades deben extraerse de la propia red

División de una tarea en sub-tareas, minimizando la comunicación entre tareas

CLUSTERING JERÁRQUICO

Procedimiento

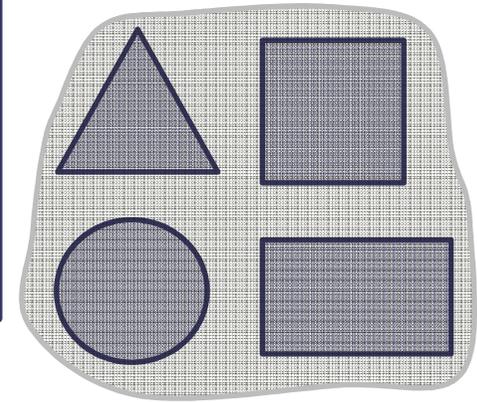
- Se calculan las similitudes entre todos los nodos (coseno, Jaccard, Euclídea, Manhattan, Hamming, etc.)
- Se comienza con todos los nodos desconectados (clustering aglomerativo)
- Se van enlazando los nodos por pares en orden de similitud decreciente
- Resultado: Componentes anidados formando un dendrograma, del que se pueden obtener distintas particiones cortando a varios niveles del árbol



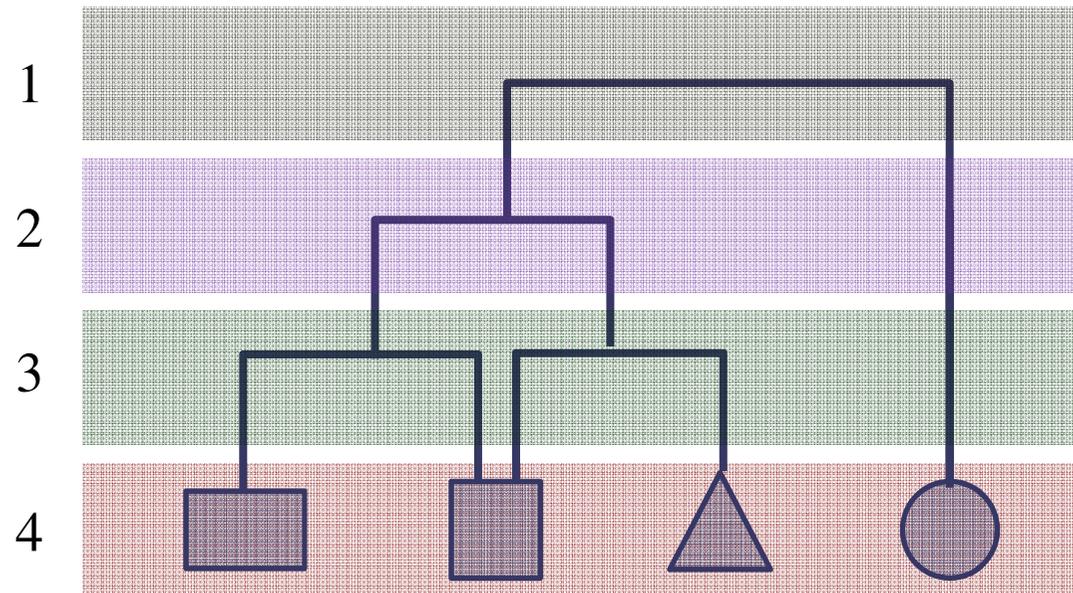
CLUSTERING JERÁRQUICO

Ejemplo ilustrativo

	Enlaces	4 Lados	Estable	Igual
1. Cuadrado	+	+	+	+
2. Rectángulo	+	+	+	--
3. Círculo	--	--	--	--
4. Triángulo	+	--	+	+

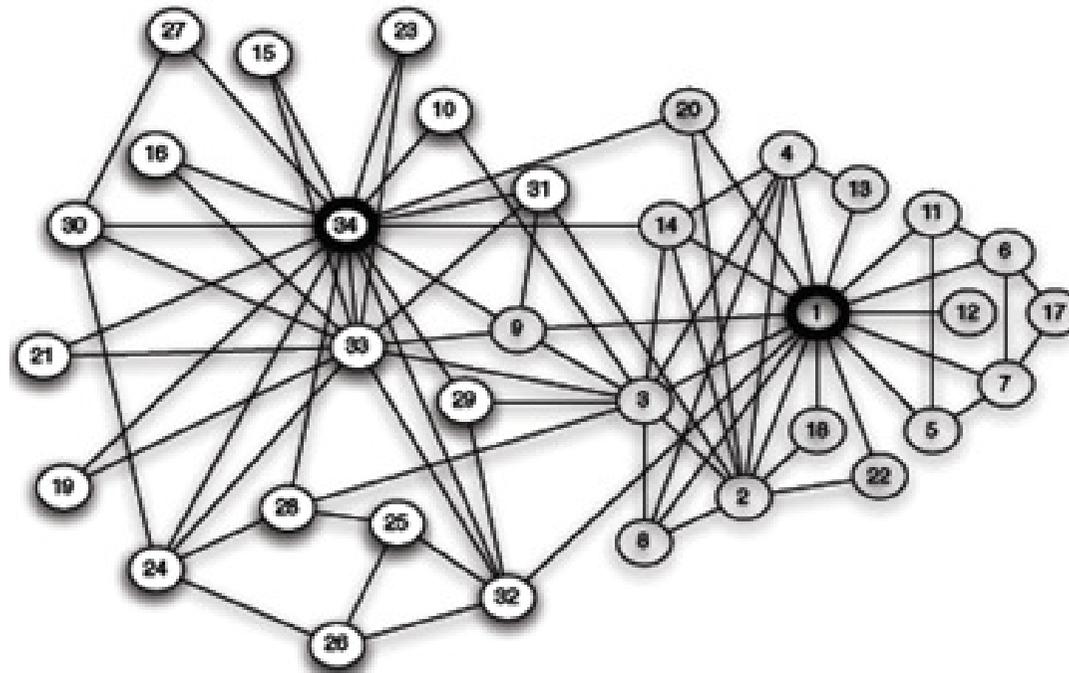


$$W_{ij} = \begin{pmatrix} 4 & 3 & 0 & 3 \\ 3 & 4 & 1 & 2 \\ 0 & 1 & 4 & 1 \\ 3 & 2 & 1 & 4 \end{pmatrix}$$



CLUSTERING JERÁRQUICO

Ejemplo: club de karate de Zachary (1)

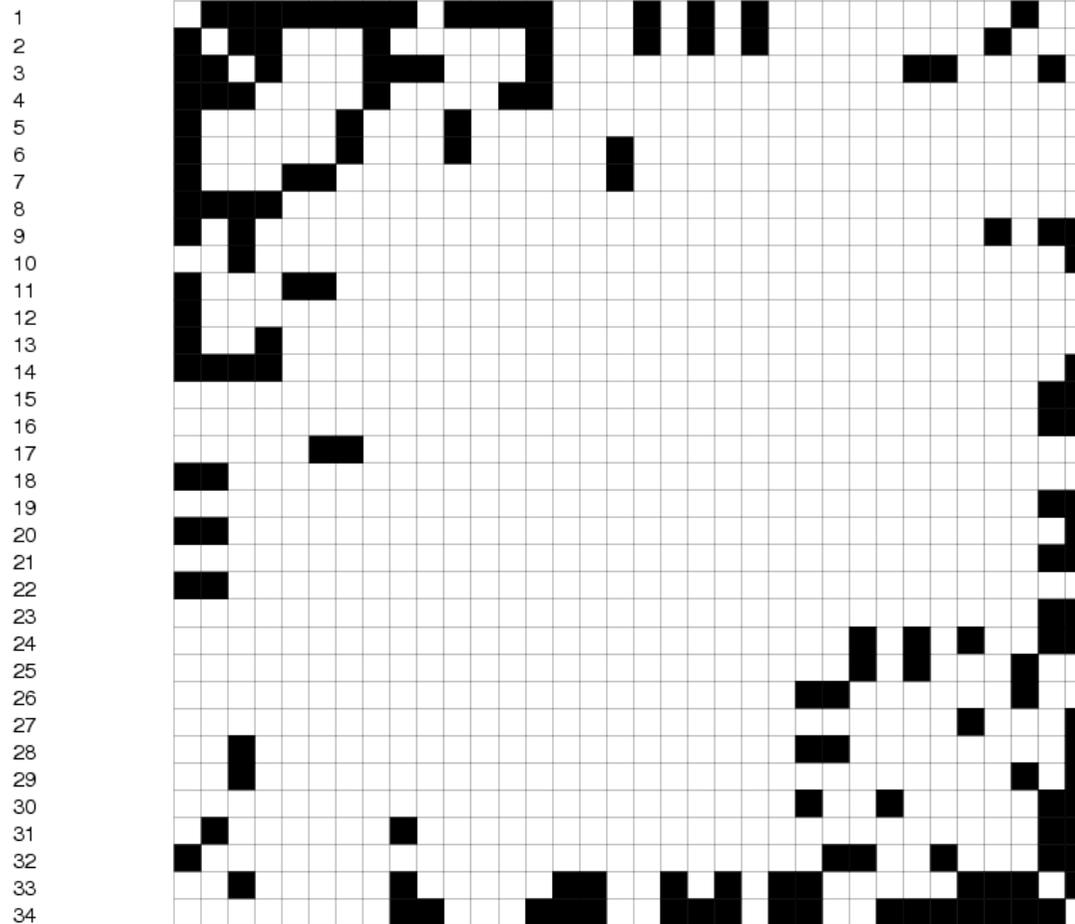


(a) *Karate club network*

Zachary. An information flow model for conflict and fission in small groups. *J. Anthropol. Res.* 33: 452-473 (1977)

CLUSTERING JERÁRQUICO

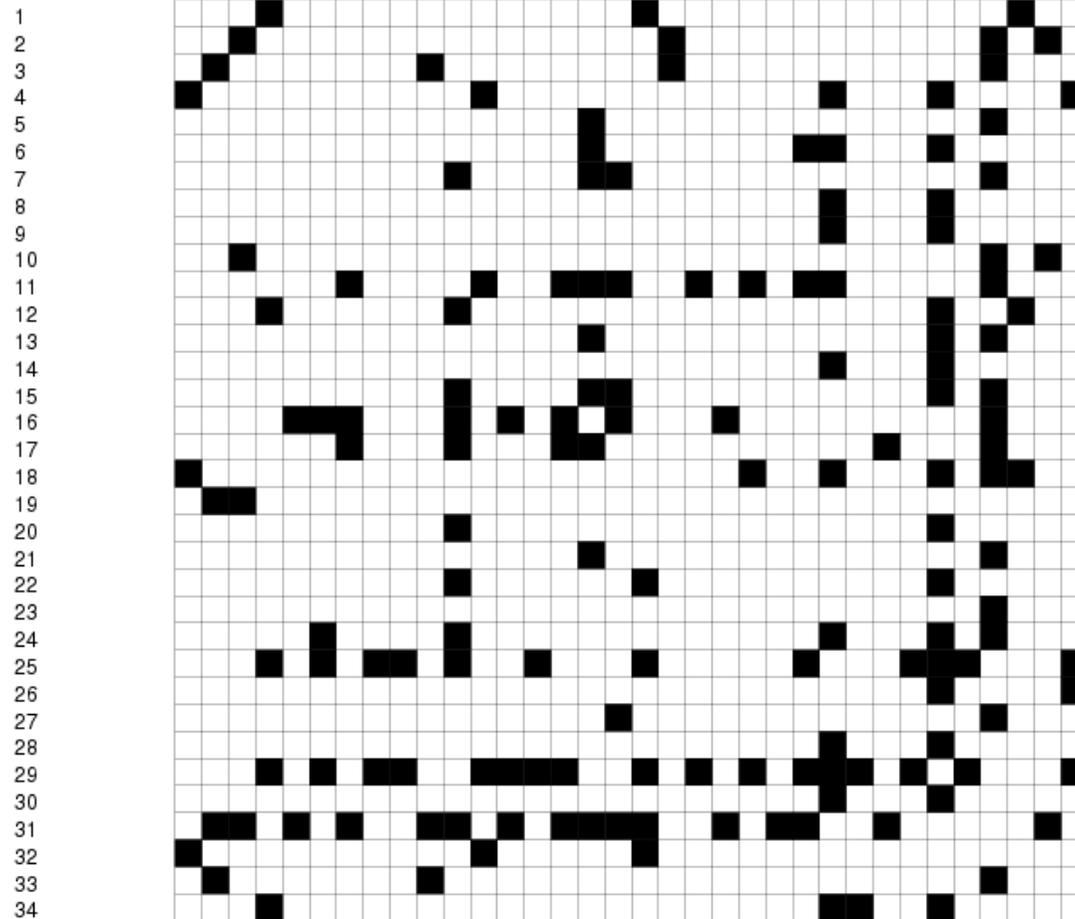
Ejemplo: club de karate de Zachary (2)



Matriz de adyacencia original

CLUSTERING JERÁRQUICO

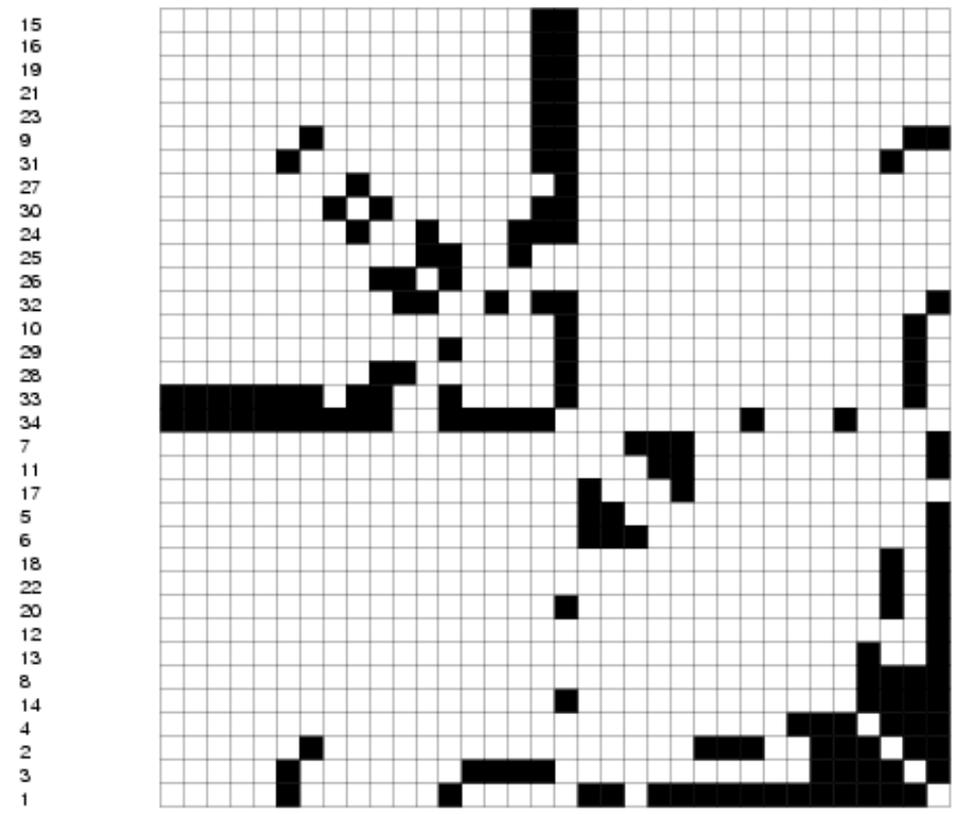
Ejemplo: club de karate de Zachary (3)



Matriz desordenada aleatoriamente

CLUSTERING JERÁRQUICO

Ejemplo: club de karate de Zachary (4)

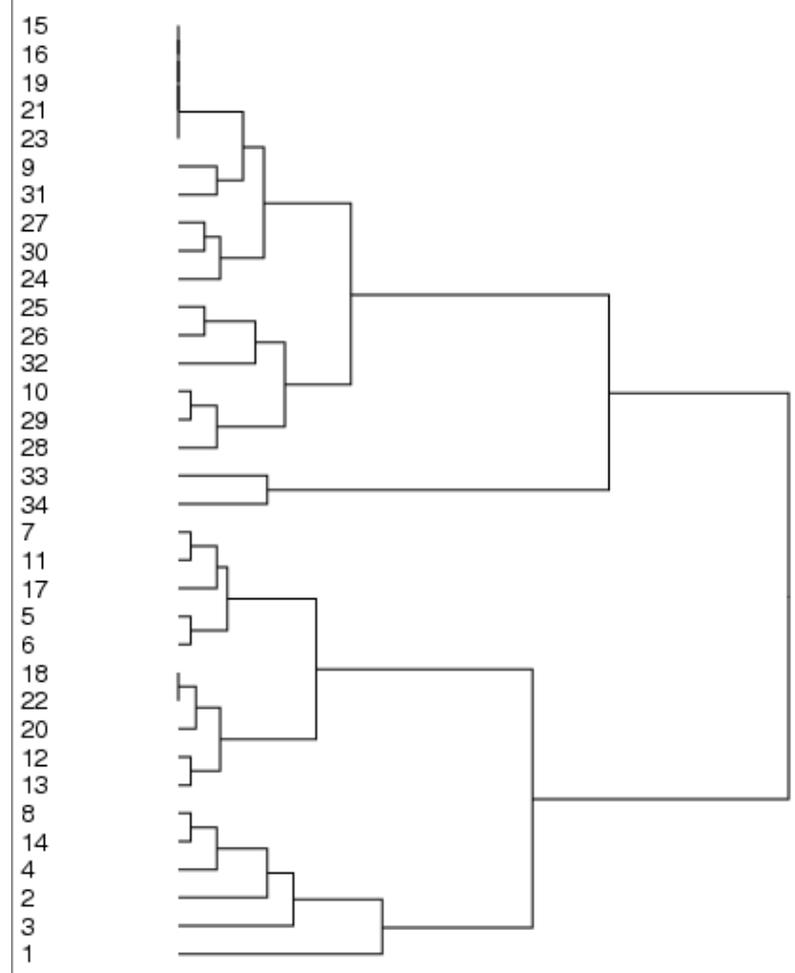


15 16 19 21 23 9 31 27 30 24 25 26 32 10 29 28 33 34 7 11 17 5 6 18 22 20 12 13 8 14 4 2 3 1

Matriz desordenada permutada

CLUSTERING JERÁRQUICO

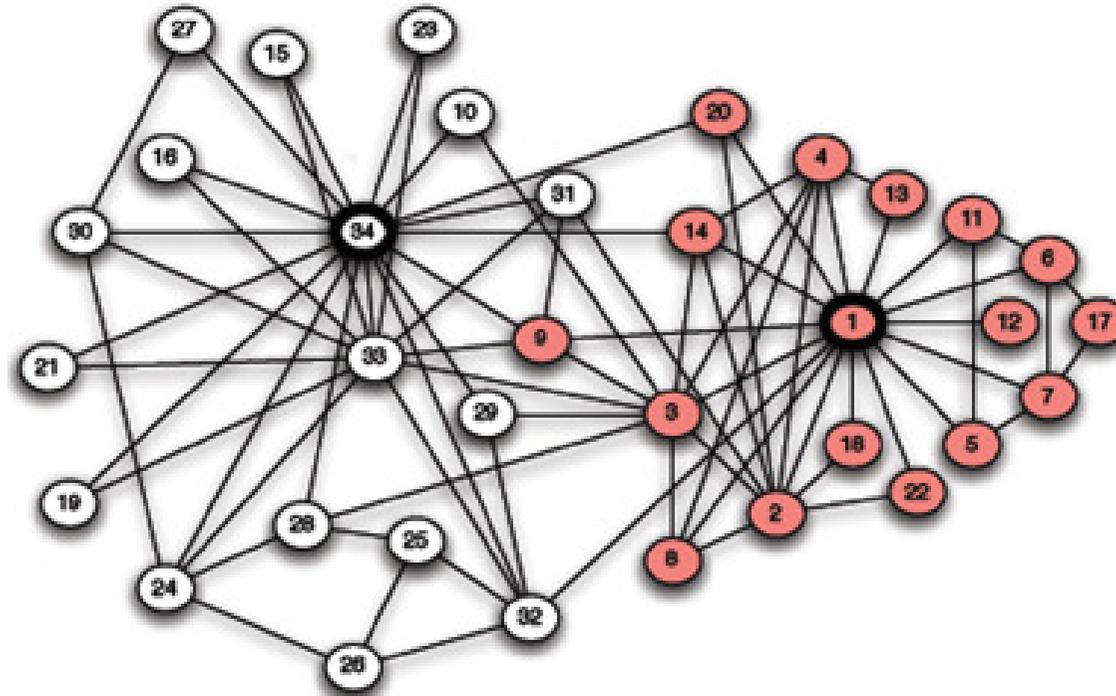
Ejemplo: club de karate de Zachary (5)



Dendrograma

CLUSTERING JERÁRQUICO

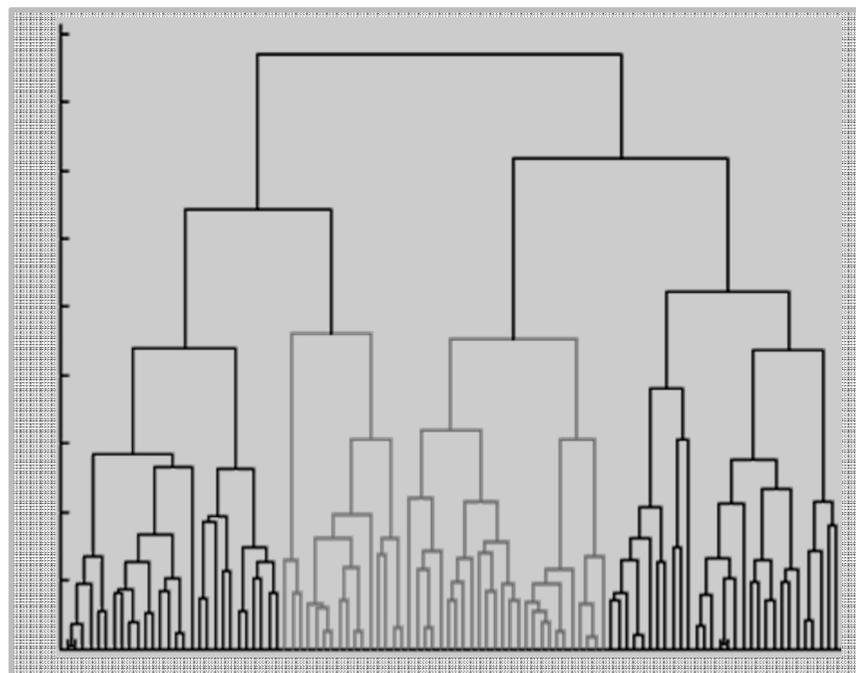
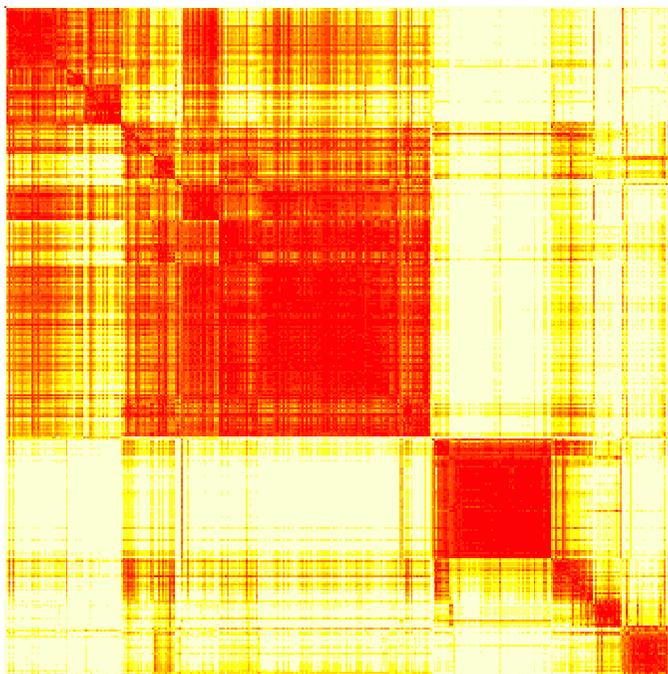
Ejemplo: club de karate de Zachary (6)



(b) *After a split into two clubs*

CLUSTERING JERÁRQUICO

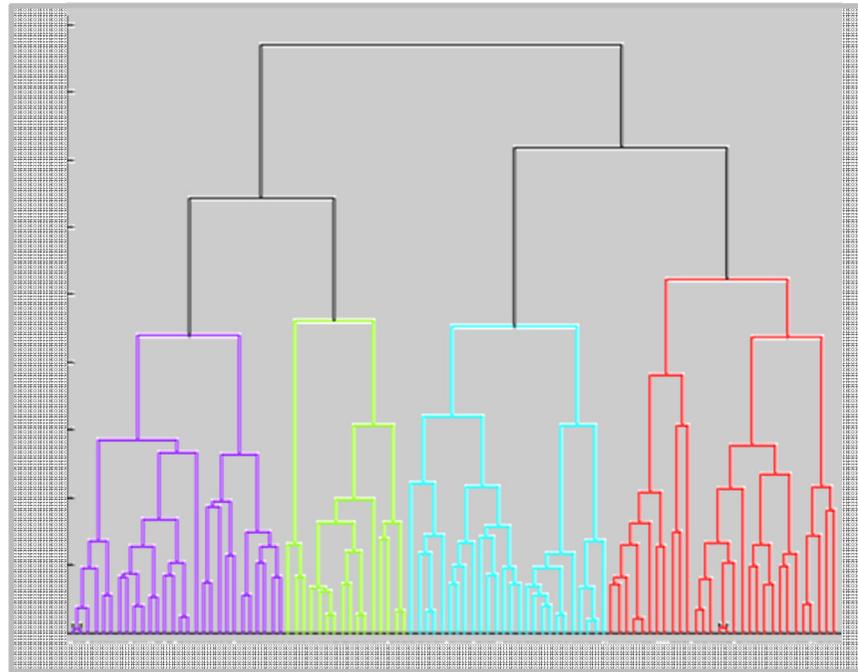
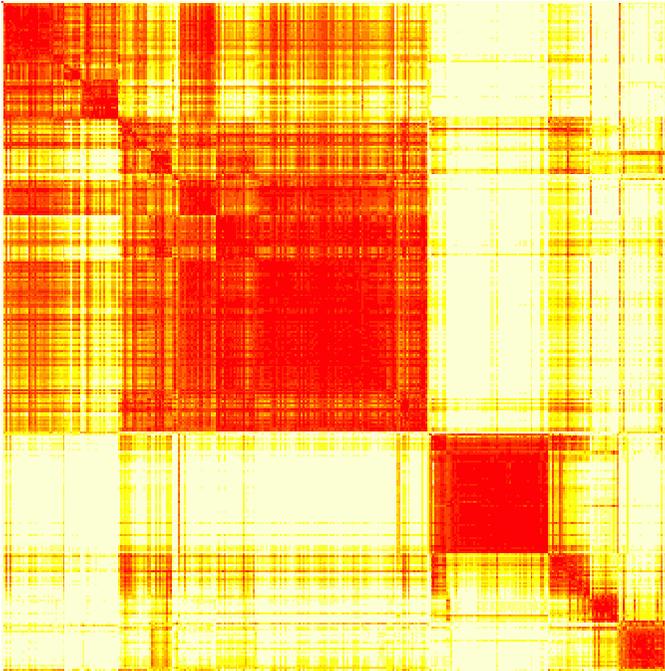
Dendrogramas (1)



Una buena estrategia para escoger la mejor partición, i.e. el mejor número de comunidades, es calcular el valor de **modularidad** para cada partición posible y seleccionar la que maximice el valor de la función

CLUSTERING JERÁRQUICO

Dendrogramas (2)



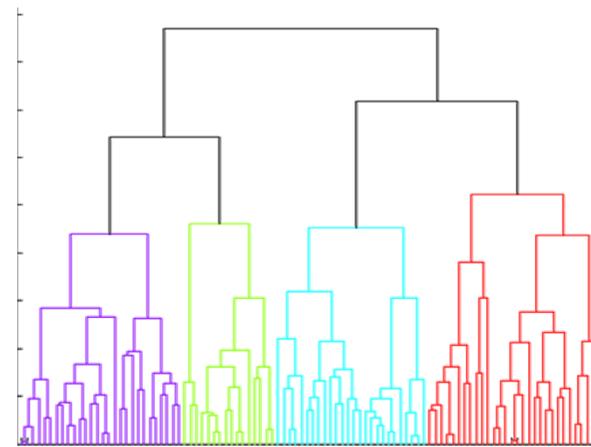
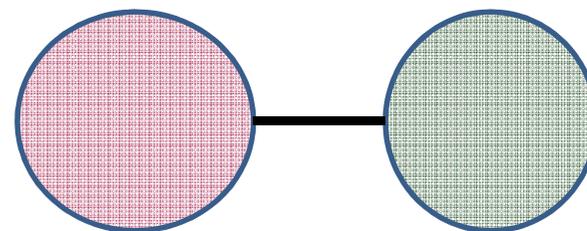
Una buena estrategia para escoger la mejor partición, i.e. el mejor número de comunidades, es calcular el valor de **modularidad** para cada partición posible y seleccionar la que maximice el valor de la función

CLUSTERING JERÁRQUICO

El método de Girvan-Newman (1)



La medida de similitud se basa en la topología de la red: **intermediación de los enlaces** (número de caminos geodésicos que pasan por el enlace en cuestión)



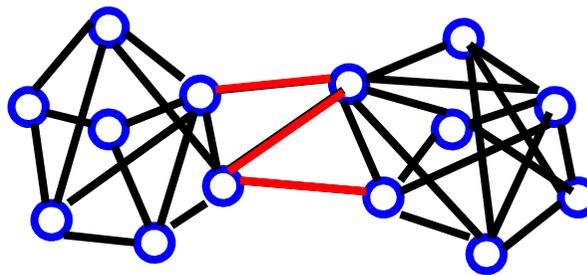
Girvan y Newman. Community structure in social and biological networks. Proc Natl Acad Sci USA 99:7821–7826 (2002)

CLUSTERING JERÁRQUICO

El método de Girvan-Newman (2)

Técnica de clustering jerárquico divisiva: divide la red original en partes conectadas más pequeñas de forma progresiva hasta que no haya más enlaces que eliminar y cada nodo representa una comunidad por sí solo o hasta un cierto umbral:

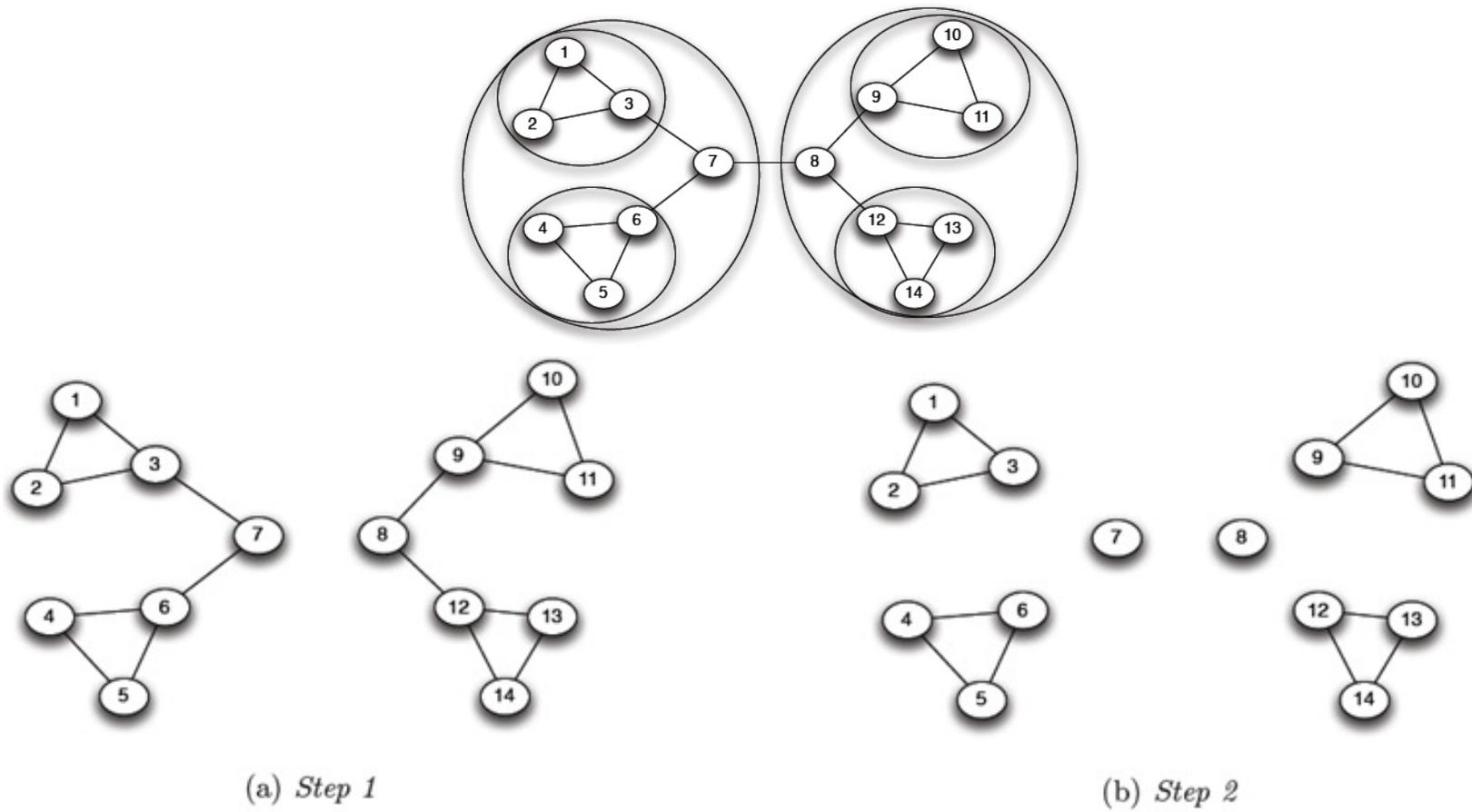
1. Calcular el valor de intermediación de todos los enlaces de la red
2. Eliminar el enlace de mayor intermediación. Este paso puede provocar que la red se divida en partes desconectadas, formando el primer nivel de nivel de regiones en la división de la red
3. Repetir los pasos anteriores hasta que no haya enlaces que eliminar en la red o hasta que el valor máximo de intermediación sea menor que un umbral



CLUSTERING JERÁRQUICO

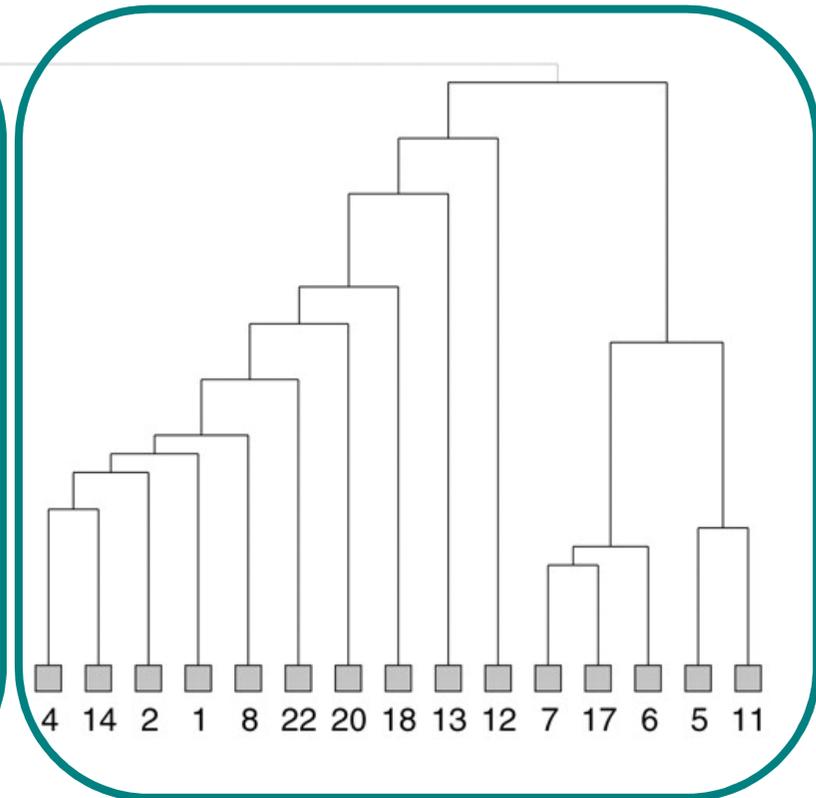
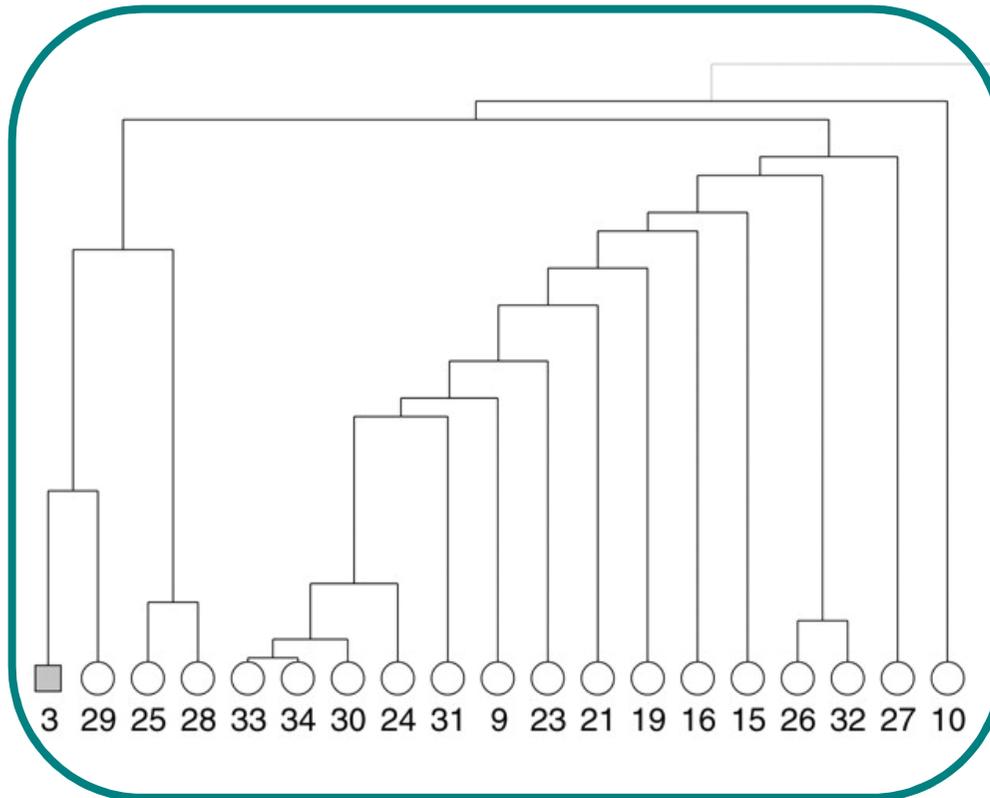
El método de Girvan-Newman (3)

Borra sucesivamente los enlaces de mayor intermediación (los puentes o los puentes locales), dividiendo la red en componentes separadas:



CLUSTERING JERÁRQUICO

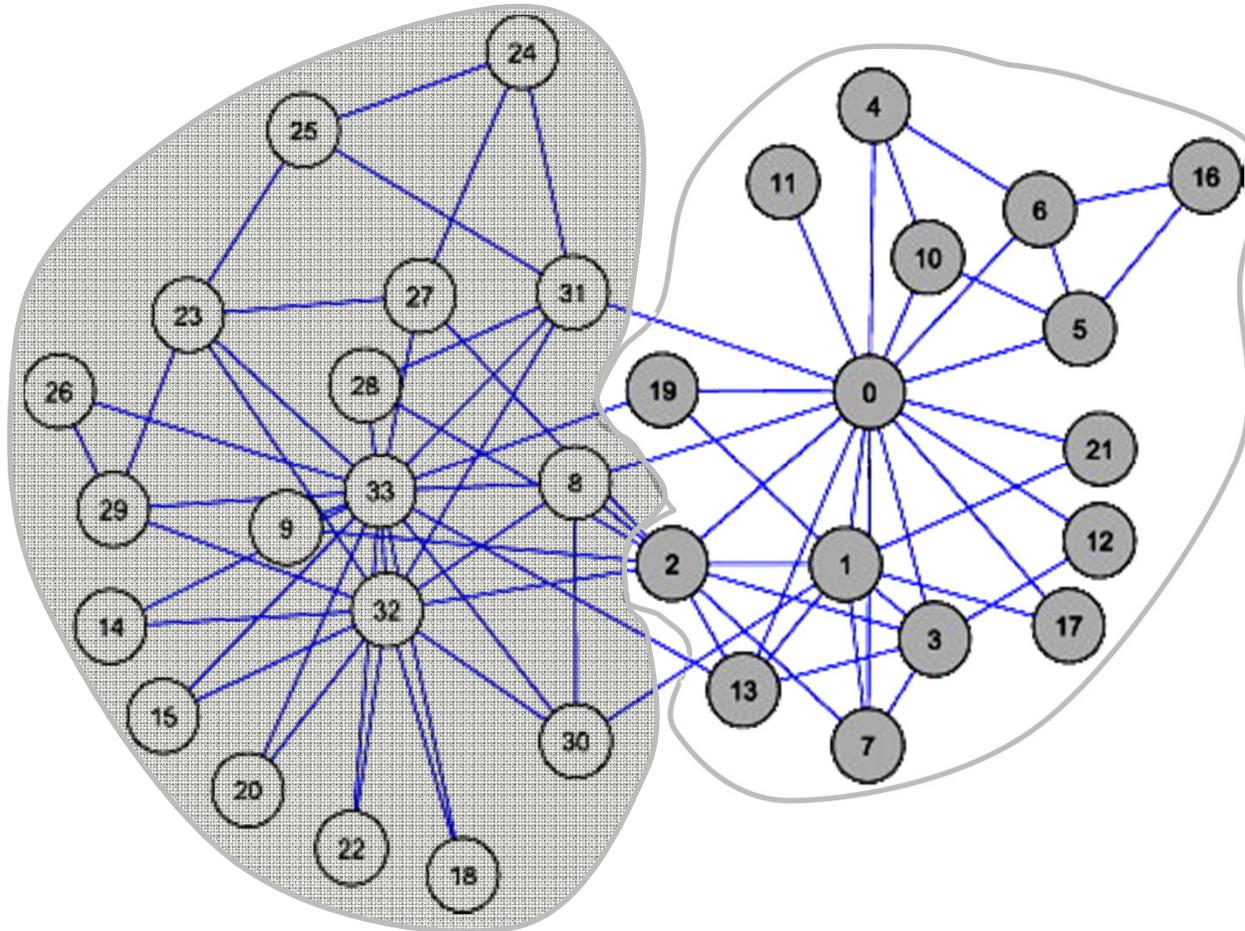
Ejemplos clásicos (1)



El club de karate de Zachary (Zachary. An information flow model for conflict and fission in small groups. *J. Anthropol. Res.* 33: 452-473 (1977))

CLUSTERING JERÁRQUICO

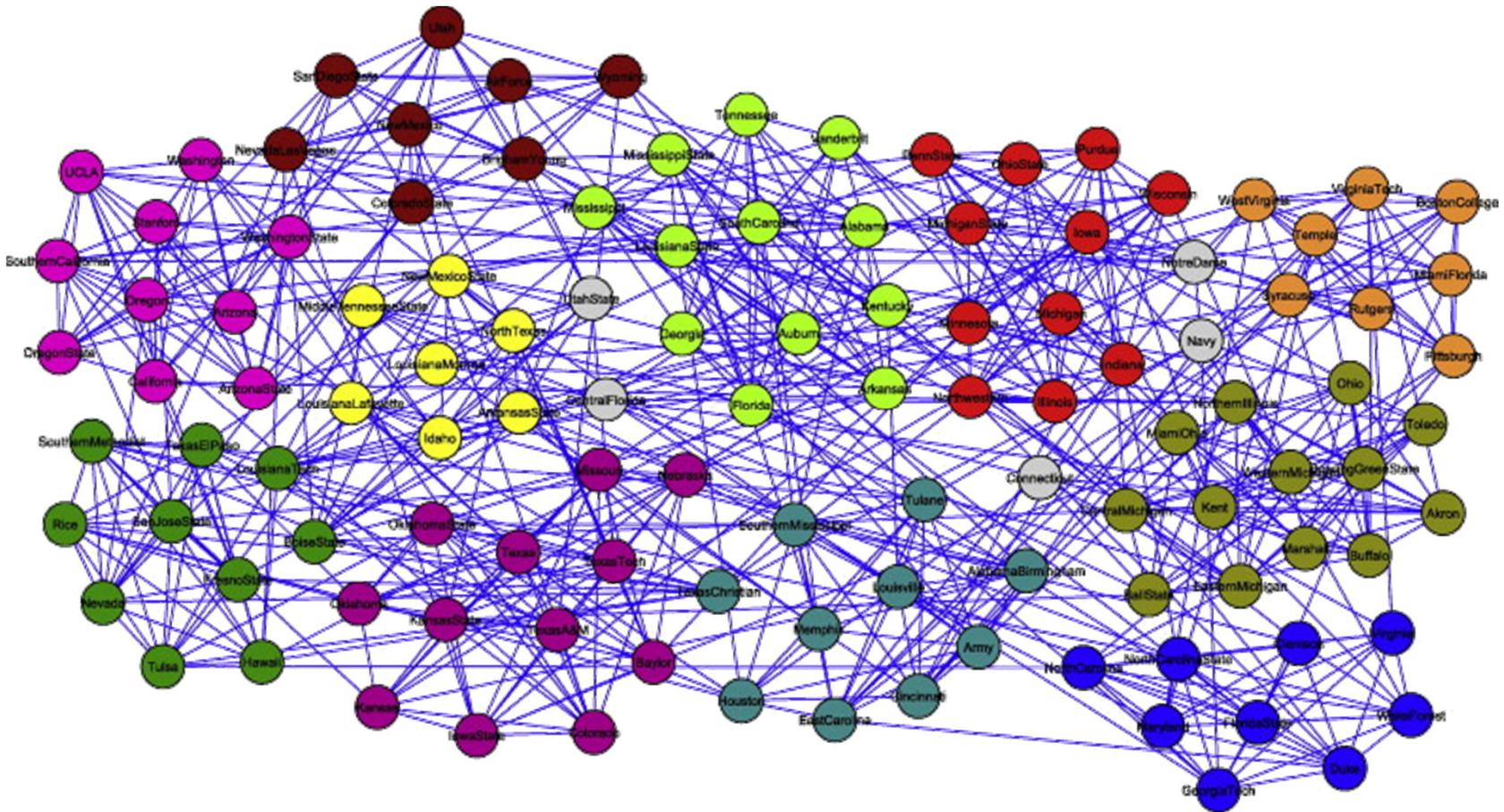
Ejemplos clásicos (2)



El club de karate de Zachary (Zachary. An information flow model for conflict and fission in small groups. *J. Anthropol. Res.* 33: 452-473 (1977))

CLUSTERING JERÁRQUICO

Ejemplos clásicos (3)



La red de fútbol americano de *colleges* (Girvan y Newman. **Community structure in social and biological networks. Proc Natl Acad Sci USA 99:7821–7826 (2002)**)

Este método de descubrimiento de comunidades es muy popular y existen un gran número de implementaciones. P.ej. la biblioteca *igraph* de *R* incorpora la función *edge.betweenness.community*

Su problema es que es necesario **recalcular la intermediación en cada paso**:

- La eliminación de un enlace puede impactar en la intermediación del resto
- Calcular los caminos geodésicos para todos los nodos es muy costoso, $O(N^3)$
- Puede ser necesario recalcularlos muchas veces
- El método no escala bien para redes con más de unos cientos de nodos, incluso en las implementaciones más rápidas

MÉTODOS BASADOS EN LA OPTIMIZACIÓN DE LA MODULARIDAD

Otro enfoque para el descubrimiento de comunidades comprende métodos heurísticos que tratan de maximizar directamente la medida de modularidad Q

Un ejemplo es el método de Blondel y otros, aplicado con éxito en redes con más de 400,000 nodos y 2 millones de enlaces, que sigue una estrategia greedy:

1. Comenzar con todos los nodos sueltos
2. Unir iterativamente los clusters que impliquen un mayor incremento ΔQ de modularidad
3. Parar cuando dicho incremento sea menor que 0

Blondel y otros. Fast unfolding of communities in large networks. J Stat Mech: Theory Exp: P10008 (2008)

Existen mejoras basadas en el uso de simulated annealing

Guimera y Amaral. Functional cartography of complex metabolic networks. Nature 433:895–900 (2005)

Ambas se utilizan en Amazon

MÉTODOS QUE CONSIDERAN EL SOLAPAMIENTO DE LAS COMUNIDADES (1)

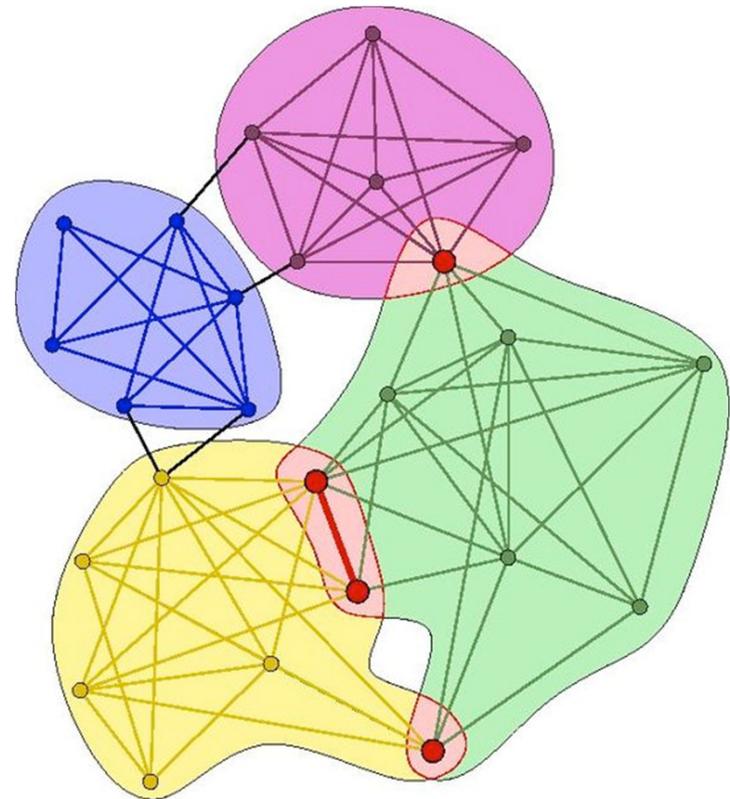
- Investigaciones recientes han demostrado que en comunidades virtuales como Orkut y Flickr los métodos de descubrimiento de comunidades no son capaces de determinar grupos de más de 100 nodos

Leskovec, Lang, Dasgupta, Mahoney. Statistical Properties of Community Structure in Large Social and Information Networks. Intl World Wide Web Conference, 2008 [Video]

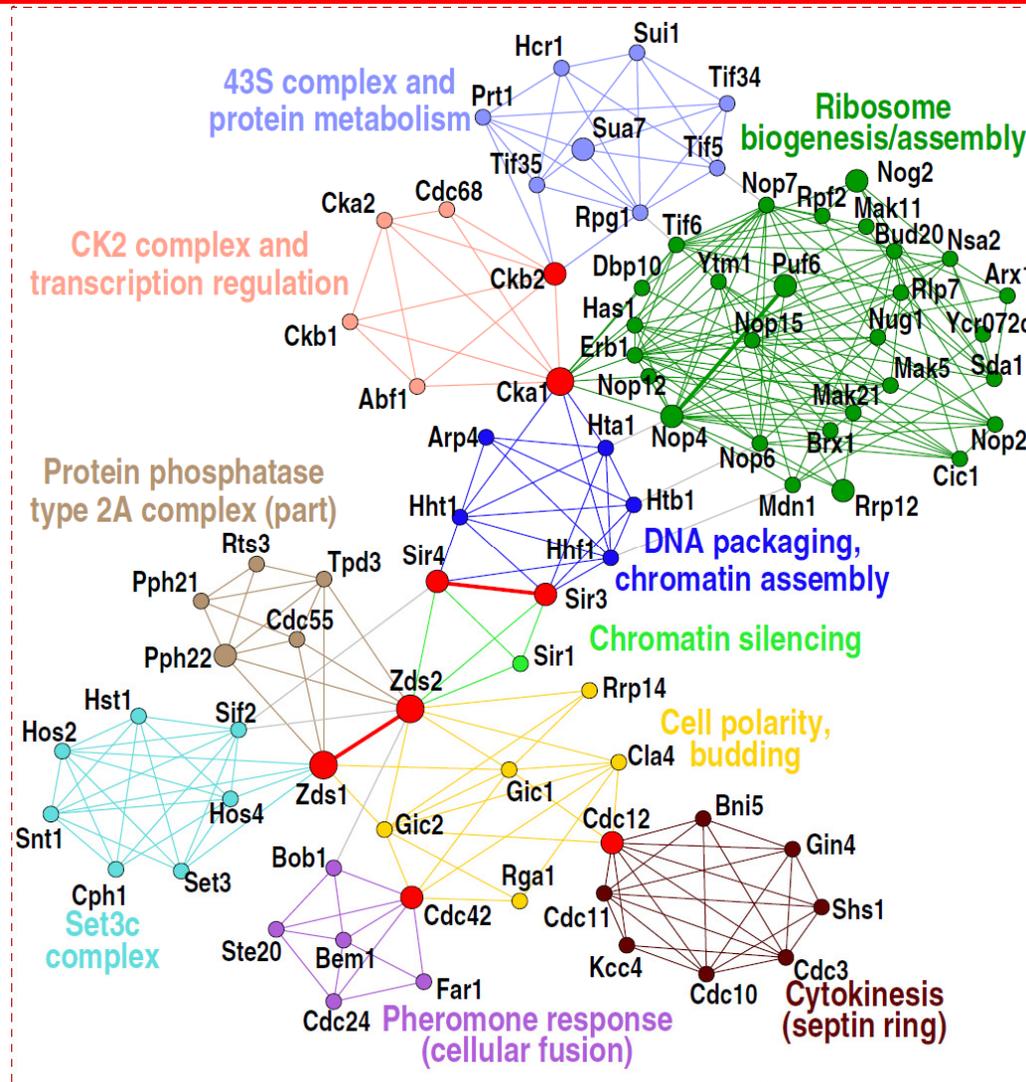
- Clique percolation method: <http://www.cfinder.org>

Palla y otros. Uncovering the overlapping community structure of complex networks in nature and society. Nature 435:814–818 (2005)

La imagen muestra comunidades identificadas por 4-cliques adyacentes. Los nodos solapados están marcados en rojo



MÉTODOS QUE CONSIDERAN EL SOLAPAMIENTO DE LAS COMUNIDADES (3)



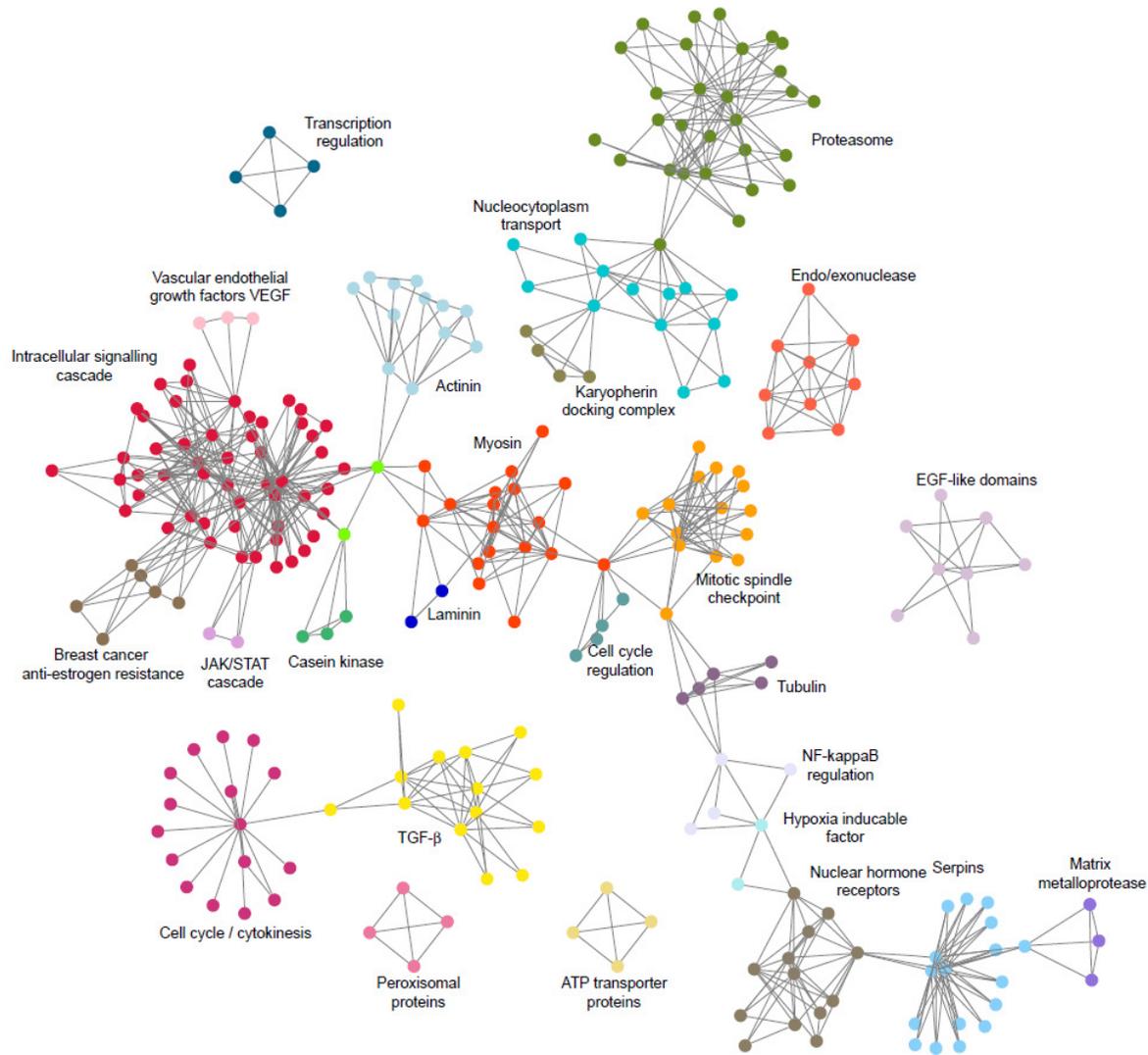
Estructura de comunidades en una red de interacción entre proteínas

Interacciones entre las proteínas de las células cancerosas de una rata

Detectadas con el método Clique Percolation

Jonsson y otros. Cluster analysis of networks generated through homology: automatic identification of important protein communities involved in cancer metastasis. *BMC Bioinf.* 7: 2 (2006)

MÉTODOS QUE CONSIDERAN EL SOLAPAMIENTO DE LAS COMUNIDADES (4)



Estructura de comunidades en una red de interacción entre proteínas

Interacciones entre las proteínas de las células cancerosas de una rata

Detectadas con el método Clique Percolation

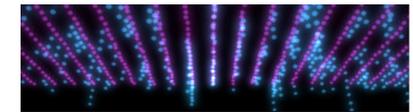
Jonsson y otros. Cluster analysis of networks generated through homology: automatic identification of important protein communities involved in cancer metastasis. BMC Bioinf. 7: 2 (2006)

Referencias y Agradecimientos

Para diseñar los materiales de este tema, he hecho uso de material desarrollado por expertos en el área disponible en Internet:

- “Network Science Interactive Book Project” del Laszlo Barabasi Lab:

<http://barabasilab.com/networksciencebook>



- Curso on-line “Social Network Analysis” de Lada Adamic, Coursera, Universidad de Michigan: <https://www.coursera.org/course/sna>

