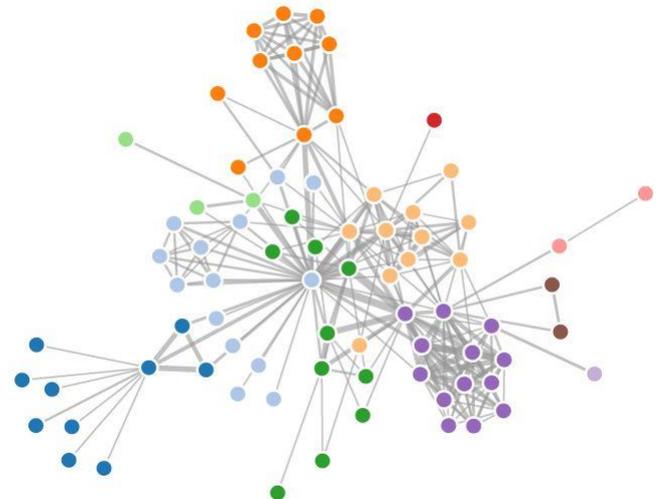
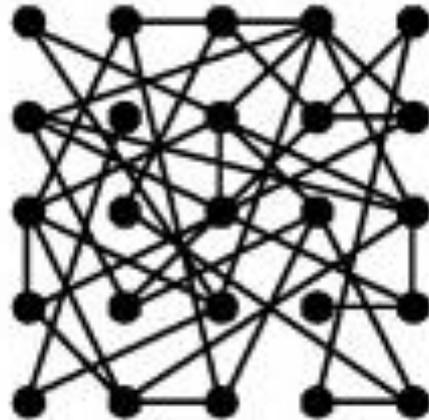
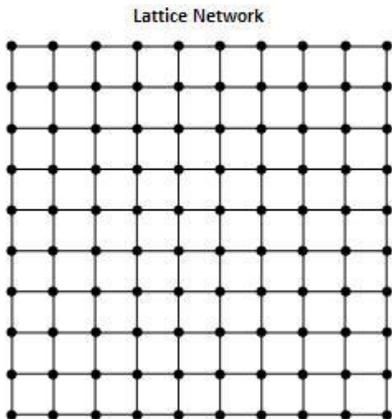


Detección de Comunidades en grafos

2 hipótesis para que funcione

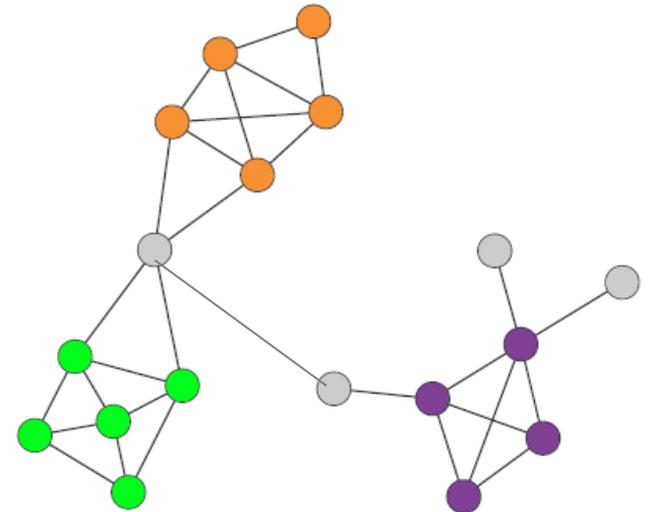
1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}

Existen realmente grupos? Existen metodologías de búsquedas, no una definición a priori de lo que buscamos. Cómo lo sabemos si no buscamos? Cómo sabemos cuales heterogeneidades son las relevantes? Cómo sabemos cuál es la **escala relevante**?



2 hipótesis para que funcione

1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
2. Qué es una comunidad? Criterio (no definición) de **Conectividad y densidad**
Una comunidad es un subgrafo **conexo, localmente denso**.
 - Desde un nodo de una comunidad puedo alcanzar cualquier otro
 - Un nodo de una comunidad se enlaza con alta probabilidad a nodos de la misma comunidad



2 hipótesis para que funcione

1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
2. Qué es una comunidad? Criterio (no definición) de **Conectividad y densidad**
Una comunidad es un subgrafo **conexo, localmente denso**.

Primeras definiciones utilizaron el concepto de **clique** o **subgrafo completo**:

subgrafo conexo de máxima densidad de enlaces



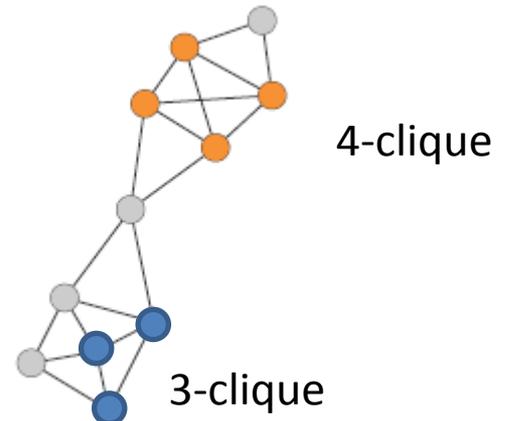
3-clique



4-clique



5-clique



El concepto de clique suele ser demasiado restrictivo en la práctica

2 hipótesis para que funcione

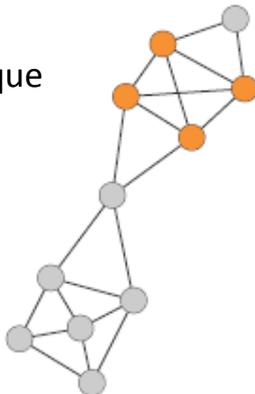
1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
2. Qué es una comunidad? Criterio (no definición) de **Conectividad y densidad**
Una comunidad es un subgrafo **conexo, localmente denso**.

Noción más laxa: C es una comunidad

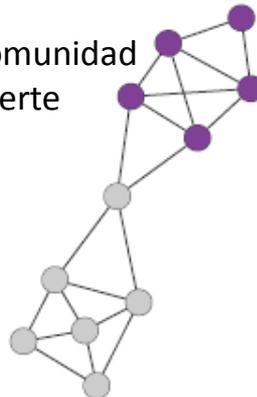
Fuerte: $k_i^{int}(C) > k_i^{ext}(C) \forall i \in C$

Débil: $\sum_{i \in C} k_i^{int}(C) > \sum_{i \in C} k_i^{ext}(C)$

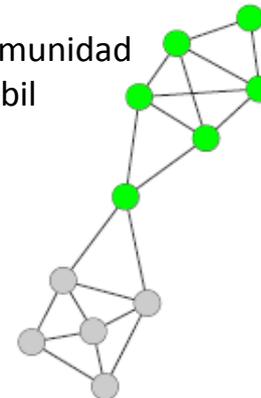
4-clique



comunidad fuerte



comunidad débil



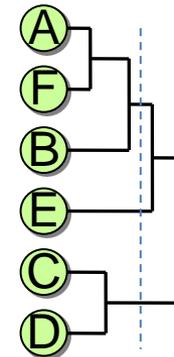
2 hipótesis para que funcione

1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
2. Qué es una comunidad? Criterio de **Conectividad y densidad**

Existen realmente grupos? Existen metodologías de búsquedas, no una definición a priori de lo que buscamos. Cómo lo sabemos si no buscamos? Cómo sabemos cuales heterogeneidades son las relevantes? Cómo sabemos cuál es la **escala relevante**?

Ya vimos algoritmo **aglomerativo**:

- I. Noción de distancia (similaridad)
- II. Agrupamiento jerárquico. Partiendo de elementos disjuntos, se adosan sucesivamente nodos y comunidades de alta similaridad
- III. Definición de grupos a partir del dendrograma



$$d_{a,c} \leq \max(d_{a,b}, d_{bc})$$

2 hipótesis para que funcione

1. **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
2. Qué es una comunidad? Criterio de **Conectividad y densidad**

Existen realmente grupos? Existen metodologías de búsquedas, no una definición a priori de lo que buscamos. Cómo lo sabemos si no buscamos? Cómo sabemos cuales heterogeneidades son las relevantes? Cómo sabemos cuál es la **escala relevante**?

Veamos ahora algoritmo **divisivo**: Girvan-Newman

Clusters à la Newman-Girvan

Idea: Partir de un cluster gigante. Ir dividiéndolo **removiendo enlaces** que conecten nodos de baja similaridad

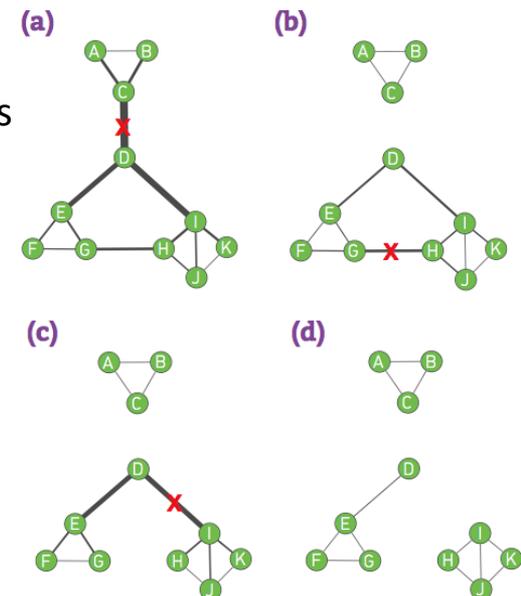
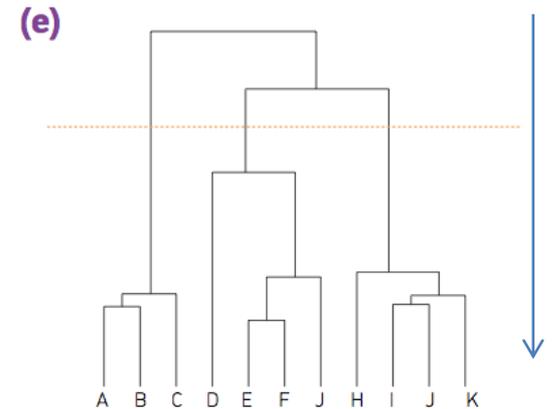
1. Definir una **centralidad de enlaces**

Intermedieatz de enlaces (*link betweenness*):

nro de caminos más cortos entre todos los pares de nodos, que atraviesen un dado enlace.

2. Agrupamiento jerárquico divisivo

- I. Computar centralidad (betweenness) de enlaces
- II. Remover el enlace de mayor centralidad
- III. Recalcular centralidad de enlaces
- IV. Repetir hasta descartar el último enlace

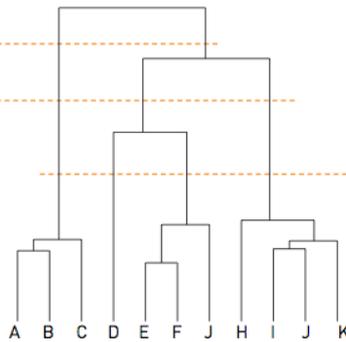


Clusters à la Newman-Girvan

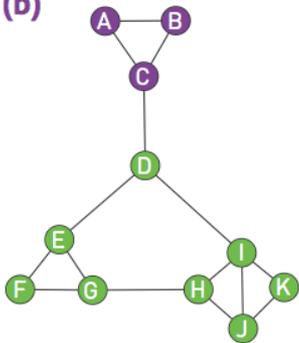
Donde **cortar** el dendrograma para definir los clusters?

Como cuantificar el **acuerdo** entre **cableado** y **partición** en grupos?

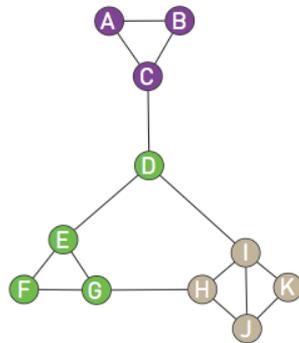
(a)



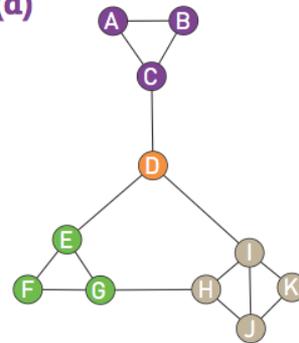
(b)



(c)



(d)

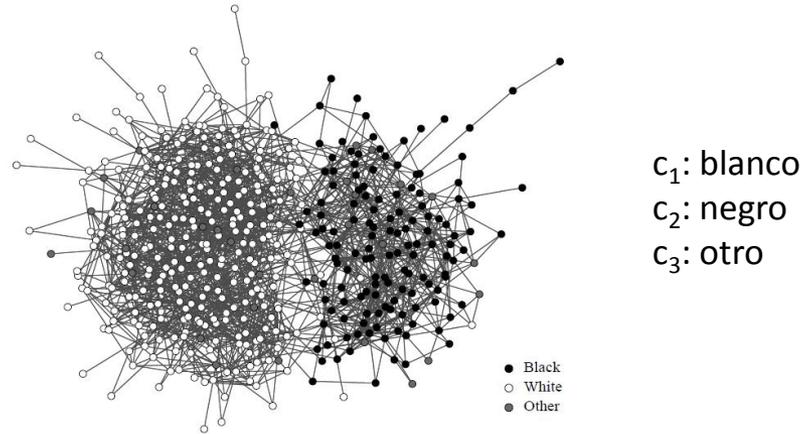


pocos enlaces internos

demasiados enlaces externos

Assortative mixing: características categóricas v.1

Supongamos que existen n_c clases diferentes para los n nodos de una red de m enlaces.
 Sea c_i la clase del nodo- i . El número de enlaces entre mismo tipo de nodos resulta:



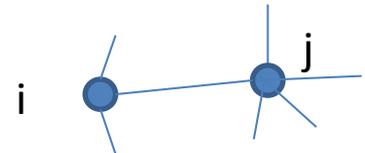
Nro de enlaces entre nodos de la misma clase

$$\sum_{\text{edges } (i,j)} \delta(c_i, c_j) = \frac{1}{2} \sum_{ij} A_{ij} \delta(c_i, c_j) \quad \delta(c_i, c_j) = 1 \text{ si } c_i = c_j$$

Red aleatoria (recableado)

$$\frac{1}{2} \sum_{ij} k_i \frac{k_j}{2m} \delta(c_i, c_j)$$

chances de que nodo- j esté al otro extremo de un enlace del nodo- i



Assortative mixing: características categóricas v.1

Supongamos que existen n_c clases diferentes para los n nodos de una red de m enlaces. Sea c_i la clase del nodo- i . El número de enlaces entre mismo tipo de nodos resulta:

Red real

$$\sum_{\text{edges } (i,j)} \delta(c_i, c_j) = \frac{1}{2} \sum_{ij} A_{ij} \delta(c_i, c_j)$$

$$\delta(c_i, c_j) = 1 \text{ si } c_i = c_j$$

Red aleatoria (recableado)

$$\frac{1}{2} \sum_{ij} \frac{k_i k_j}{2m} \delta(c_i, c_j).$$

Consideramos la diferencia: $\frac{1}{2} \sum_{ij} A_{ij} \delta(c_i, c_j) - \frac{1}{2} \sum_{ij} \frac{k_i k_j}{2m} \delta(c_i, c_j)$

modularidad

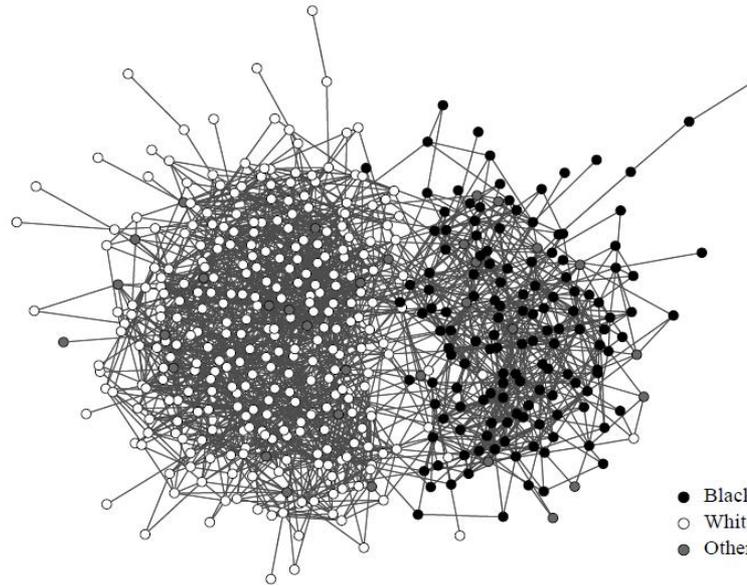
$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

- $0 < Q < 1$ si hay más enlaces entre vertices del mismo tipo que los esperados por azar
- $Q < 0$ si hay menos enlaces entre vertices del mismo tipo que los esperados por azar

DEJA VÍ

Assortative mixing: características categóricas

DEJA VÍU



- Black
- White
- Other

fracción enlaces adyacentes a nodos de categoría r

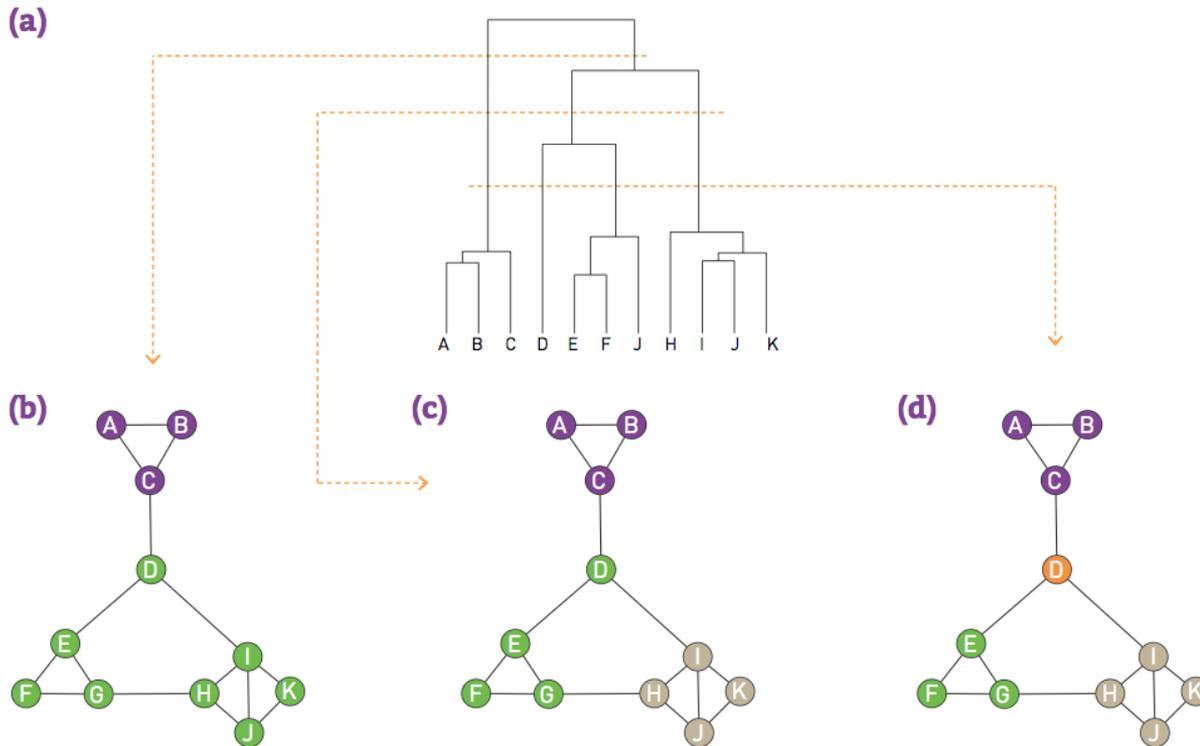
$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j).$$

$$Q = \sum_r (e_{rr} - a_r^2)$$

fracción enlaces entre nodos de categoría r

Clusters à la Newman-Girvan

Donde **cortar** el dendrograma para definir los clusters?



pocos enlaces internos

demasiados enlaces externos

Como cuantificar el **acuerdo** entre **cableado** y **partición** en grupos?

fracción enlaces adyacentes a nodos de cluster r

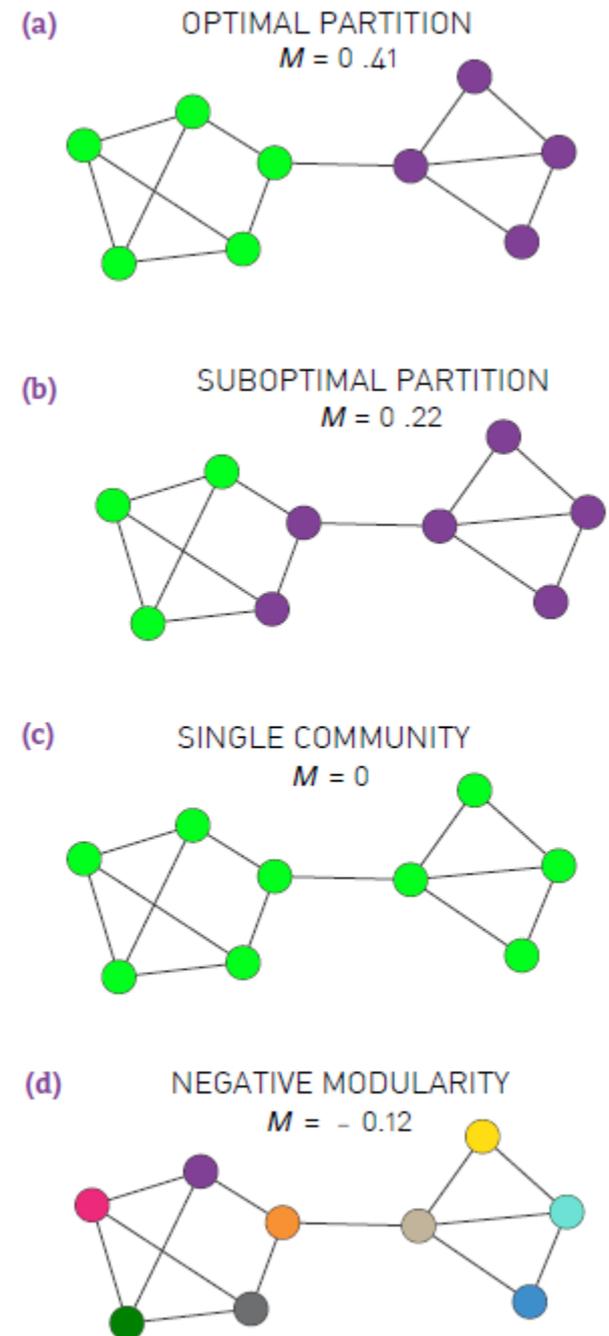
$$Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$$

fracción enlaces entre nodos de cluster r

Propiedades de la modularidad

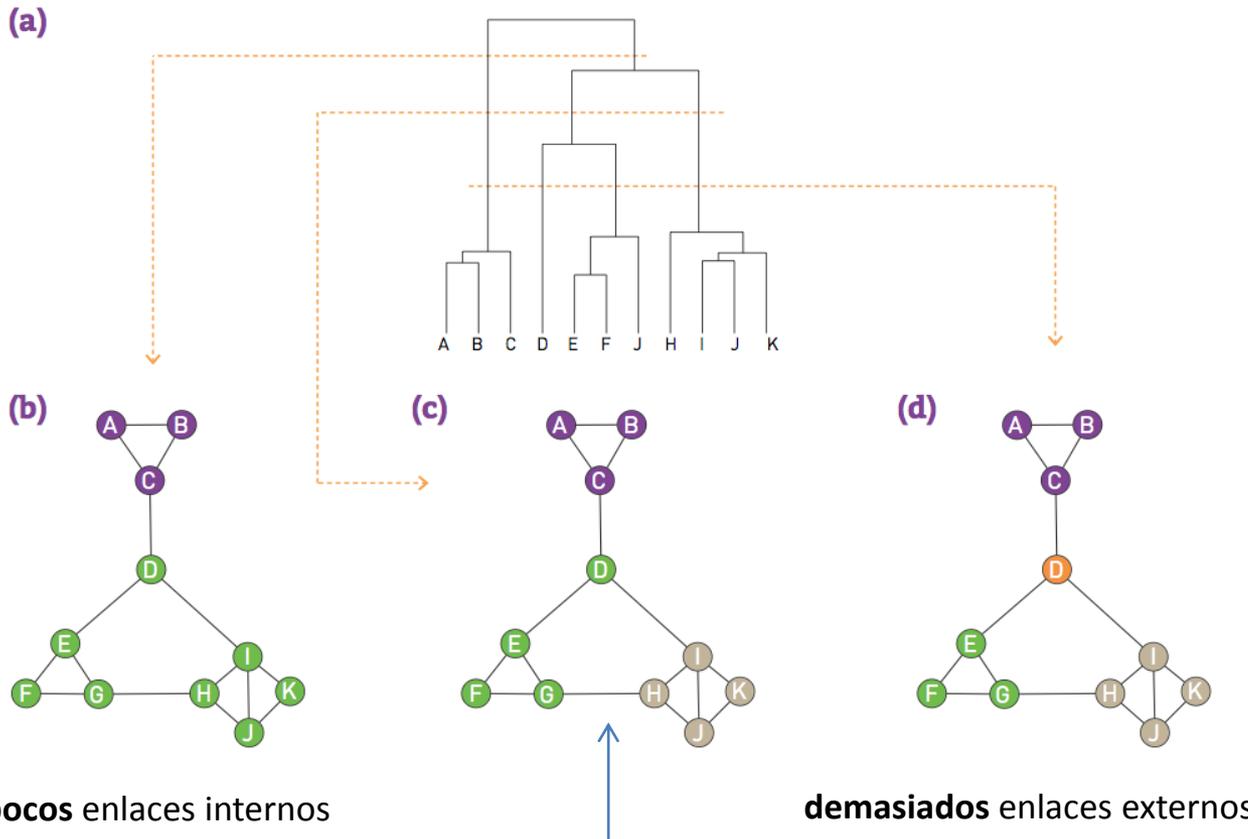
$$Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$$

- Valores altos de modularidad implican mejor valor de asortatividad entre pertenencia a grupos y cableado de la red.
- Una partición de un único cluster tendrá $Q=0$ (los dos términos son idénticos)
- Si cada nodo pertenece a una comunidad distinta $L_r=0$ y $Q<0$



Clusters à la Newman-Girvan

Donde **cortar** el dendrograma para definir los clusters?



Como cuantificar el **acuerdo** entre **cableado** y **partición** en grupos?

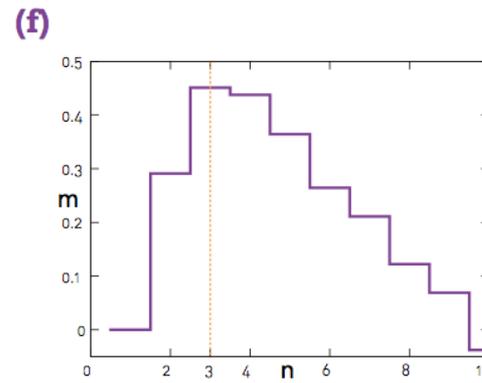
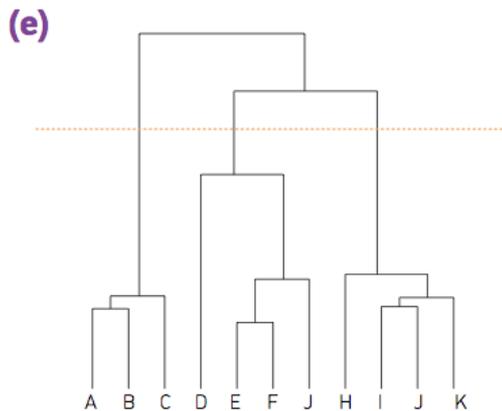
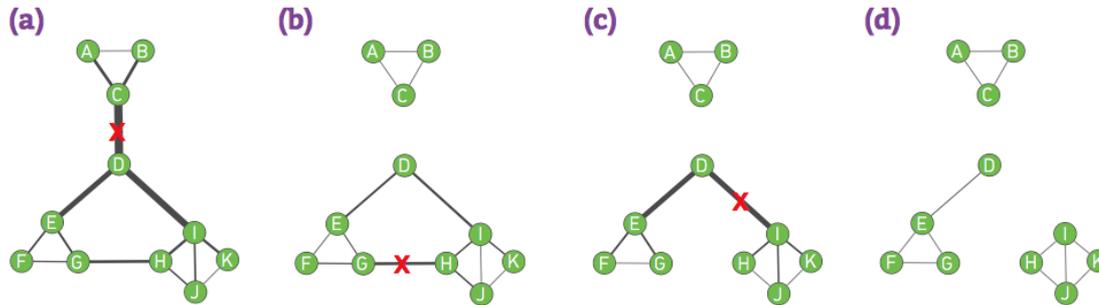
fracción enlaces adyacentes a nodos de cluster r

$$Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$$

fracción enlaces entre nodos de cluster r

Maxima modularidad

Entonces...cual partición?



$$Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$$

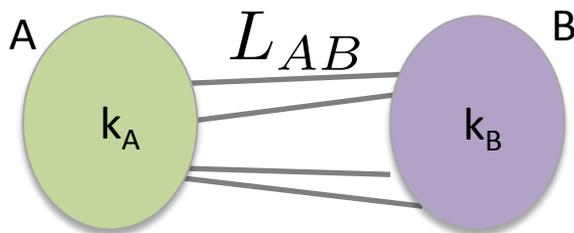
Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)

Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)

Supongamos que tenemos dos grupos y analizamos juntarlos $Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$



k_x : grado total de grupo X
 L_A : nro enlaces en A
 L_B : nro enlaces en B
 L_{AB} : nro enlaces entre A y B

$$\Delta Q = \left[\frac{L_{AB} + L_A + L_B}{L} - \left(\frac{k_{AB}}{2L} \right)^2 \right] - \left[\frac{L_A}{L} - \left(\frac{k_A}{2L} \right)^2 + \frac{L_B}{L} - \left(\frac{k_B}{2L} \right)^2 \right]$$

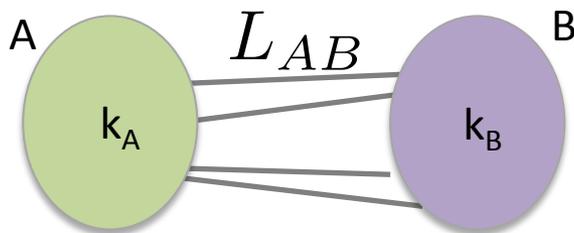
$$\Delta Q = \left[\frac{L_{AB}}{L} - \frac{k_A k_B}{2L^2} \right]$$

Si $\Delta Q > 0$ se debería promover la fusion

Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)

Supongamos que tenemos dos grupos y analizamos juntarlos $Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$



k_x : grado total de grupo X
 L_A : nro enlaces en A
 L_B : nro enlaces en B
 L_{AB} : nro enlaces entre A y B

$$\Delta Q = \frac{L_{AB}}{L} - \frac{k_A k_B}{2L^2} \quad (9.57) \text{ Barabasi}$$

No sólo depende del cableado de A y B, sino del número de enlaces, L, de la red completa !

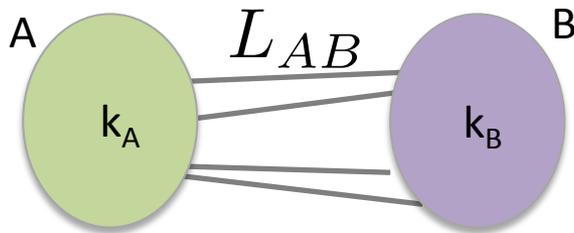
Si $L > \sqrt{\frac{k_A k_B}{2}} \Rightarrow \Delta Q > 0$

y maximizar Q implicará unir ambos clusters **SIEMPRE !!!**
 (aunque $L_{AB}=1$)

Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)

Supongamos que tenemos dos grupos y analizamos juntarlos $Q = \sum_r \left(\frac{L_r}{L} - \left(\frac{k_r}{2L} \right)^2 \right)$



k_x : grado total de grupo X
 L_A : nro enlaces en A
 L_B : nro enlaces en B
 L_{AB} : nro enlaces entre A y B

$$\Delta Q = \frac{L_{AB}}{L} - \frac{k_A k_B}{2L^2} \quad (9.57) \text{ Barabasi}$$

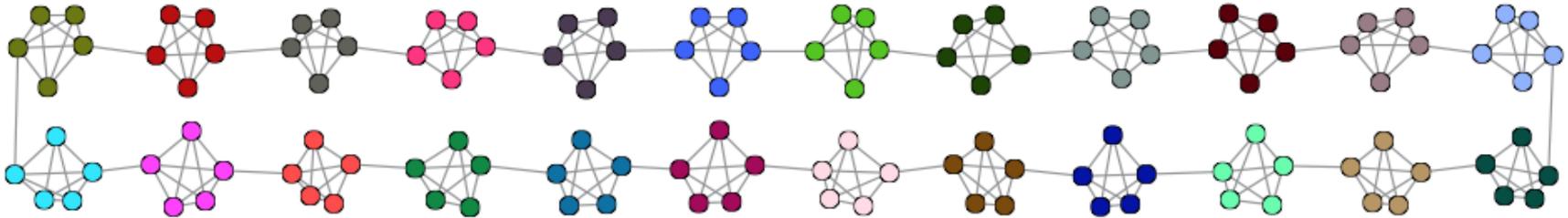
Si $L > \sqrt{\frac{k_A k_B}{2}} \Rightarrow \Delta Q > 0$
 y maximizar Q implicará unir ambos clusters SIEMPRE (aunque $L_{AB}=1$)

Si asumimos $k_A \sim k_B \equiv \kappa \Rightarrow \kappa < \sqrt{2L}$

Algoritmos que maximicen Q **no podrán identificar** comunidades de **tamaño** menor a κ (!)
 Se suele hacer zoom...

Limite de resolucion...ejemplo

Partición *natural* para una red conformada por 24 5-cliques



Pero...en esta red, unir dos cliques es negocio!

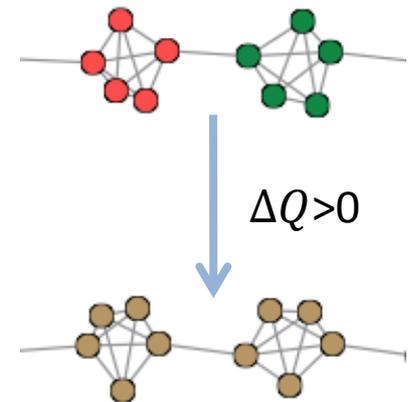
$$\Delta Q = \frac{L_{AB}}{L} - \frac{k_A k_B}{2L^2}$$

$$= \frac{L_{AB}}{L} - \frac{1}{L} \frac{k_A k_B}{2L} = \frac{1}{L} \left(L_{AB} - \frac{k_A k_B}{2L} \right) > \frac{1}{L} (L_{AB} - 0.75)$$

$$L = 10 * 24 + 24 = 264$$

$$k_A = k_B = 20$$

$$\frac{k_A k_B}{2L} = \frac{400}{528} < 1$$

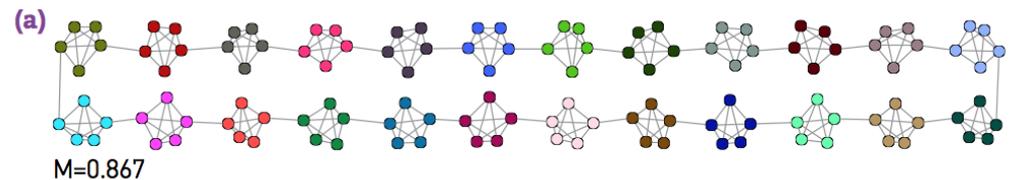


$\Delta Q > 0$ si $L_{AB} \geq 1$

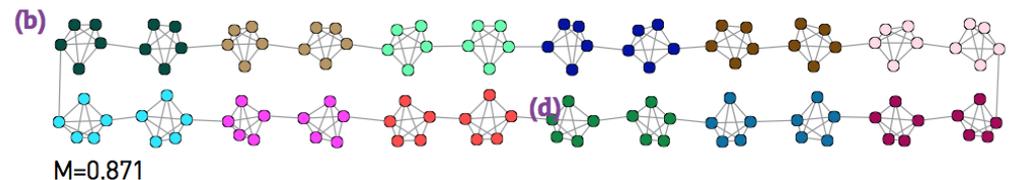
Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)

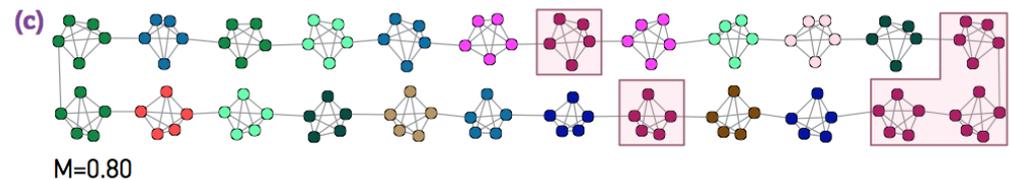
Partición *natural*



Partición *optima*

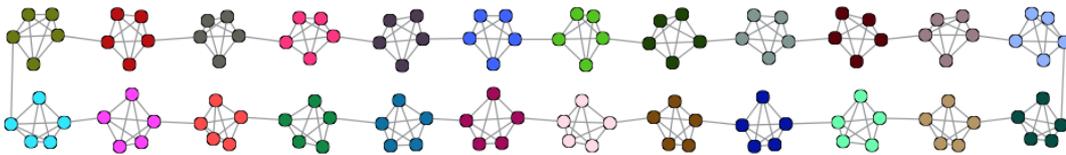


Partición *aleatoria*

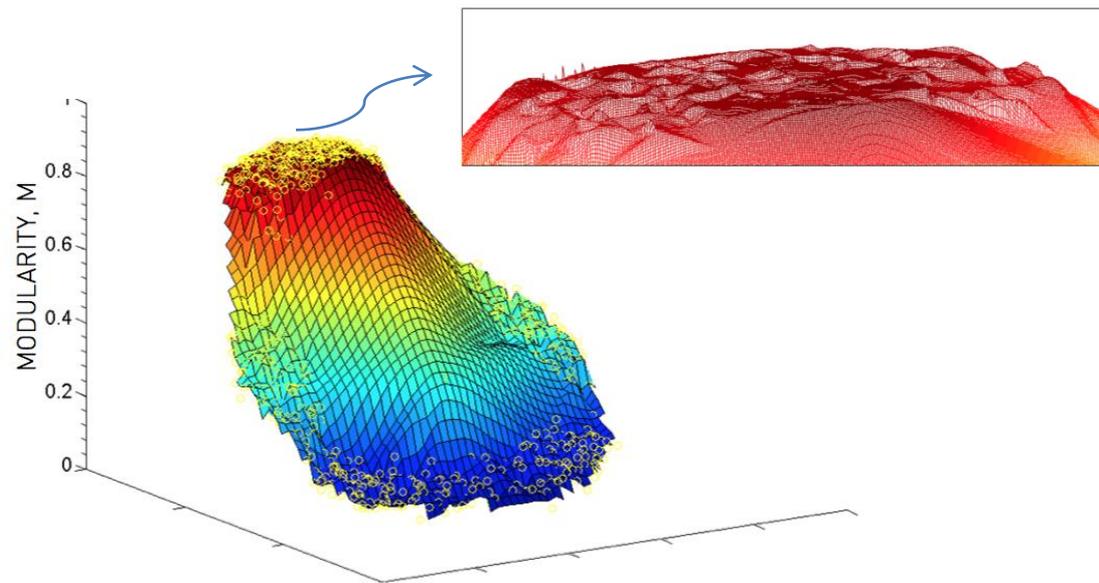


Limitaciones de la Modularidad

- Límite de resolución
- Máximo global, no muy diferente de máximos locales (*spin-glass landscape*)



Modularidad estimada para 997 particiones. No se observa un máximo claro para Q . Particiones muy diferentes en su naturaleza, podrían presentar valores comparables de Q



4 hipótesis para que funcione

- **Suponemos** que existe una estructura en comunidades embebida en la red, i.e. en A_{ij}
- Criterio de **Conectividad y densidad**. Una comunidad se corresponde con un subgrafo conexo localmente denso
- Redes aleatorias carecen de una estructura de comunidades
- La modularidad de una partición sobre una red permite identificar particiones óptimas en el sentido de mixing asortativo entre pertenencia a comunidades y conexas (entender las limitaciones de esto!)

