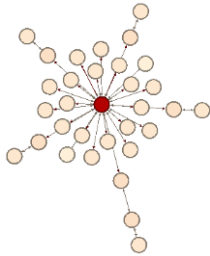


Jorge Fábrega

Redes Sociales



1

Objetivos de la clase

- Introducir aplicaciones estadísticas clave en teoría de redes que permiten modelar **formación y evolución de vínculos**.
- Mostrar cómo distintos modelos (ERGMs, Ising, EGA) capturan **regularidades estructurales, dinámicas y latentes**.
- Conectar intuiciones sociales con fundamentos matemáticos

2

¿Por qué estudiar estadísticas en redes?

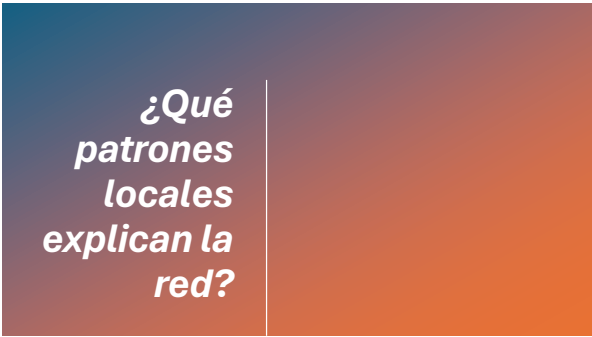
- Las redes sociales no son aleatorias: muestran **regularidades estructurales** (reciprocidad, homofilia, cierre triádico).
- Los grafos empíricos son solo una **muestra parcial**: necesitamos inferir y generalizar.
- Nos interesa explicar **cómo emergen configuraciones colectivas** a partir de decisiones locales.
- Ante eso: La estadística nos permite:
 - Comparar la red observada con un **mundo contrafactual** (qué pasaría "por azar").
 - Identificar **mecanismos sociales ocultos** detrás de los enlaces.

3

Tres enfoques principales

- **Estructurales**
 - Pregunta: *¿Qué patrones locales explican la red?*
 - Ejemplos: ERGMs, block models.
 - Enfoque: regularidades en triángulos, reciprocidad, centralidad.
- **Dinámicos**
 - Pregunta: *¿Cómo cambian los enlaces en el tiempo?*
 - Ejemplos: TERGMs, preferential attachment.
 - Enfoque: evolución de triadas, coevolución de atributos y vínculos.
- **Latentes**
 - Pregunta: *¿Qué estructuras invisibles generan la red?*
 - Ejemplos: Ising, Expected Graph Models (EGM).
 - Enfoque: detectar dimensiones ocultas (ideología, rasgos psicológicos, etc.).

4



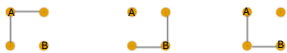
5



Exponential
Random Graph
Models



6



¿cuál es la probabilidad de que A sea amigo de B?

¿qué combinación de lazos de amistad en toda la clase es más probable, dadas las características de los nodos y de la red en su conjunto?

7

ERGMs



- Es el "equivalente" en redes a un general lineal model
- Expected Random Graph Model (ERGM) es una familia de modelos que permite testear hipótesis sobre qué explica la formación/características de una red.
 - Por ejemplo: Si para una red G , tu hipótesis implica que:
 - La cantidad de conexiones importa → testear densidad
 - La reciprocidad importa → testear relaciones no dirigidas
 - La formación de comunidades importa → testear triadas
 - Los atributos de los nodos importa → testear su distribución, homofilia,...
 - Etc...

8

La intuición



- "Es como" una regresión logística, pero para redes.
- Básicamente, buscamos probabilidad de que exista un link entre los nodos i y j , $P(a_{ij}=1)$ dado el resto de la estructura del grafo en las características que consideras que explican la red, G . Es decir: $P(a_{ij}=1|G)$
- El output de un ERGM es un odd-ratio para cada una de esas características del siguiente modo:
$$\log \left(\frac{P(a_{ij}=1|G)}{P(a_{ij}=0|G)} \right)$$
- Por ejemplo, si tu modelo tiene una sola característica, digamos, el número de links y el valor que obtuviste fue: -1.19, la probabilidad es 0,23
 - ¿Por qué? Porque es un odd-ratio, debes transformarlo: $\exp(-1.19)/[1+\exp(-1.19)]$

9

Ahora formalmente



- Sea $G=(V,E)$ una red. Sea $Y_{ij} = Y_{ji}$ una variable binaria que adquiere valor 1 si $e_{ij} \in E$.
- Entonces $Y = [Y_{ij}]$ es la *matriz aleatoria adyacente* de G . Tal que $y = [y_{ij}]$ es una realización particular de dicha matriz.
- Un ERGM es la distribución de probabilidades conjunta de los elementos de $Y = y$:
$$P_{\theta}(Y = y) = \left(\frac{1}{k}\right) \exp\left\{\sum_H \theta_H g_H(y)\right\}$$
- Donde: H es una “**configuración**” o set de posibles edges entre los vértices o nodos de G
- $g_H(y) = \prod_{y_{ij} \in H} y_{ij}$, por lo tanto es 1 si la configuración H ocurre y 0 si no.
- k es una constante de normalización: $k = k(\theta) = \exp\{\sum_H \theta_H g_H(y)\}$

10

¿Por qué una exponencial?



- Porque las distribuciones exponenciales tienen una característica muy atractiva: NO TIENEN MEMORIA.
- (una breve explicación aquí:
<https://lilc.stat.purdue.edu/2014/41600/notes/prob3204.pdf>)
- Entonces, podemos estudiar las probabilidades condicionales de cada link independiente de los otros, “como si” fuesen una cadena de eventos independientes.

11

Teorema de Hammersley–Clifford (1971)



- Toda distribución de probabilidad puede representarse como una cadena de Markov si cumple con ciertos requisitos en la distribución de probabilidades.
- Las distribuciones exponenciales cumplen con esas propiedades.
- Implicancia: toda familia de modelos de redes puede ser expresado (potencialmente) como un ERGM si se cuenta con las estadísticas fundamentales (S_1, S_2, \dots) de dicha familia de modelos

$$Pr(g) = \frac{e^{\theta_1 S_1(g) + \theta_2 S_2(g) + \dots + \theta_k S_k(g)}}{\sum_{g' \in G} e^{\theta_1 S_1(g') + \theta_2 S_2(g') + \dots + \theta_k S_k(g')}}.$$

12

MCMC: Markov Chain Monte Carlo



- En redes muy pequeñas ($n < 8$) es posible estimar todas las configuraciones para un set de hipótesis/configuraciones (véase paquete *ergm* en R)
- En redes más grandes deben hacerse aproximaciones estadísticas
- MCMC
 - Markov Chain: Sólo el estado inmediatamente anterior importa (separar estimación en bloques)
 - Monte Carlo: Aleatorización para sacar estimaciones de parámetros
 - Juntos: Dados los valores de estadísticos de la red, simulaciones sucesivas para buscar la distribución que mejor explica las estadísticas observadas

13

Pasos para la estimación vía ERGMs

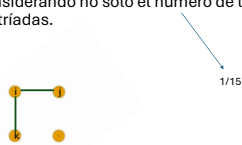


- Paso 1: Cada link es tratado como una variable aleatoria
- Paso 2: Un conjunto de hipótesis se postula para explicar formación de links
 - Ejemplo: links aleatorios, links por cercanía, por homofilia, triadic closure, etc...
- Paso 3: Cada hipótesis es una configuración que se suma a otras posibles configuraciones.
 - Ejemplo: El link ij puede ser por aleatoridad + por cercanía + por homofilia, ...
- Paso 4: Las configuraciones, aunque locales, son homogéneas en toda la red (reducción de parámetros)
- Paso 5: Estimación e interpretación de los parámetros del modelo

14

Ejemplo

- Supongamos que queremos estudiar simultáneamente la probabilidad del grafo g considerando no sólo el número de links sino también el número de tríadas.



15

Ejemplo

Aplicado a una red de Erdős-Renyi (en la que todos los links son independientes unos de otros, i.e. $\theta_T=0$):

$$Pr(g) = p^m(1-p)^{\binom{n}{2}-m} = \frac{p^m}{(1-p)^m} * (1-p)^{\binom{n}{2}}$$
$$Pr(g) = e^{\log\left(\frac{p}{1-p}\right)*m + \text{constante}}$$

$Pr(g) = e^{\theta_m m(g) + \text{constante}}$ sobre todos los posibles grafos g'

16

Ejemplo

• Entonces, sea $m(g)$ el número de links y $T(g)$ el número de triadas de la red g :
La probabilidad del grafo g depende de:

$$\begin{aligned} &\theta_m m(g) + \theta_T T(g) \\ &\quad \downarrow \\ &e^{\theta_m m(g) + \theta_T T(g)} \\ &\quad \downarrow \\ &Pr(g) = \frac{e^{\theta_m m(g) + \theta_T T(g)}}{\sum_{g' \in G} e^{\theta_m m(g') + \theta_T T(g')}} \quad \leftarrow \text{ERGM} \end{aligned}$$

17

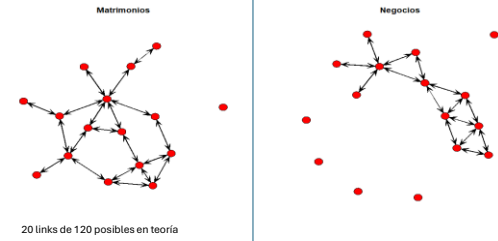
¿Cómo se hace en la práctica?

En R se hace con el paquete statnet (compilado de varios paquetes)
Referencias: <https://statnet.org/trac>

- Algoritmo:
1. Se estiman $\theta_1, \theta_2, \theta_3 \dots$ (varios métodos, maximum likelihood)
 2. Se simulan redes con propiedades $\theta_1, \theta_2, \theta_3 \dots$
 3. Se actualizan los valores de $\theta_1, \theta_2, \theta_3 \dots$ a partir de esas redes simuladas
 4. Se repite 2... hasta lograr convergencia

18

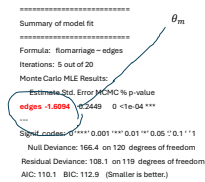
Ejemplo: Medici



19

L 479

Medici – ergm en base al número de links



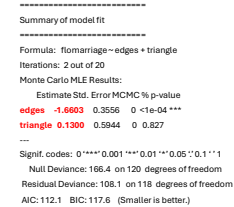
Odd ratio = -1,609 * 1 (cambio de 0 a 1)

$$\Pr(l_{ij} = 1) = \frac{e^{-1,609}}{1 + e^{-1,609}} = 0,1667$$

... que es justamente lo que esperábamos porque:
20 links reales /120 posibles = 0,1667

20

Medici – ergm en base al número de links + triádas



Odd ratio = -1,609 * 1 + 0,13*T
Si T = 0 → odd ratio: -1,609
Si T = 1 → odd ratio: -1,609+0,13=-1,47

$$\Pr(l_{ij} = 1|T = 1) = 0,185$$

... y sube en aprox.2 pts por triáda

21

Summary of model fit

Formula: flomarriage ~ edges + triangle + nodecov("wealth")

Iterations: 2 out of 20

Monte Carlo MLE Results:

Estimate Std. Error MCMC % p-value

edges -2.580037 0.536265 0 <1e-04 ***

triangle -0.116286 0.620783 0 0.8517

nodecov.wealth 0.010686 0.005018 0 0.0321 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null Deviance: 166.4 on 120 degrees of freedom

Residual Deviance: 103.1 on 117 degrees of freedom

AIC: 109.1 BIC: 117.4 (Smaller is better)

Como *wealth* son atributos de nodos, se consideran ambas riquezas al ver la probabilidad de cada link *ij*.

Wealth:

[1] 10 36 55 44 20 32 8 42 103 48 49 3

27 10 146 48

Ej: link con una triada entre familias con riqueza 10 y 146:

-2,58 - 0,11 + 0,01*(10+146) = -1,13

En prob: 0,24

Este modelo sugiere mayor probabilidad de un link a mayor riqueza de ambas familias

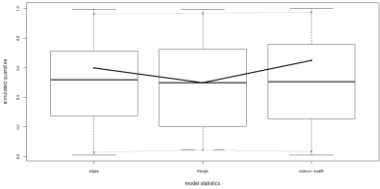
22

Goodness of fit en los Medici

¿el modelo cómo se ajusta a los datos?

A partir del modelo se simulan redes y se calculan las mismas variables del modelo

Goodness-of-fit diagnostics



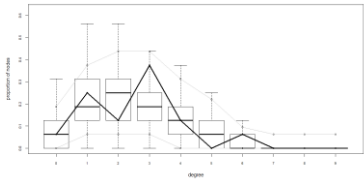
23

Goodness of fit en los Medici

Y también otras que son descriptores básicos de la red, pero no estaban en el modelo:

Distribución de grado

Goodness-of-fit diagnostics



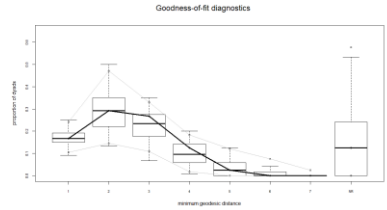
24

8

Goodness of fit en los Medici

Y también otras que son descriptores básicos de la red, pero no estaban en el modelo:

Geodésicas

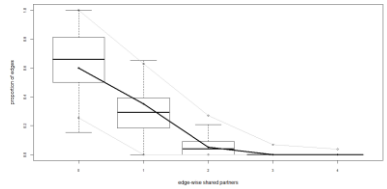


25

Goodness of fit en los Medici

Y también otras que son descriptores básicos de la red, pero no estaban en el modelo:

Vecinos compartidos



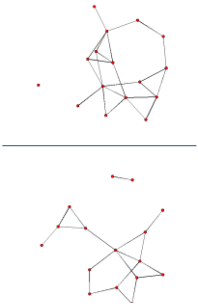
26

Simulaciones de redes

• Para poder hacer GOF fue necesario hacer simulaciones.

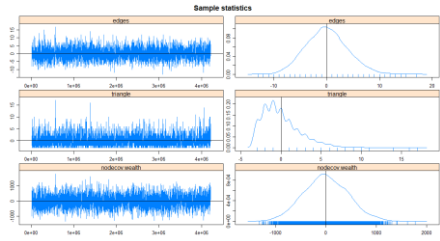
Stored network statistics:

	edges	triangle	nodecov.	wealth
[1,]	23	5		2374
[2,]	20	4		1961
[3,]	22	4		2585
[4,]	19	0		2092
[5,]	18	0		1494
[6,]	20	5		2378
[7,]	21	2		2308
[8,]	19	2		1661
...				



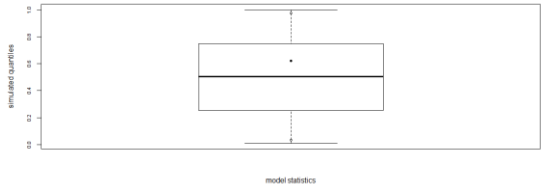
27

Diagnóstico



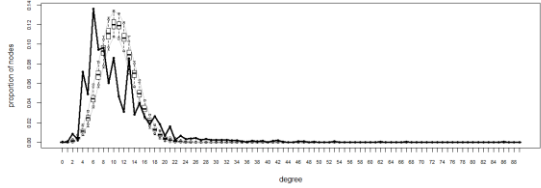
28

Un grafo aleatorio ajustado a los datos de interlocking (personas)



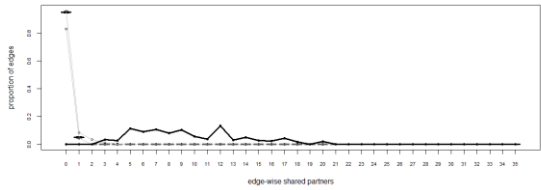
29

Un grafo aleatorio ajustado a los datos de interlocking (personas)



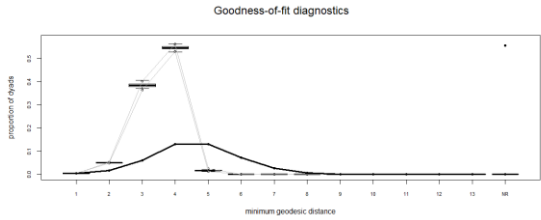
30

Un grafo aleatorio ajustado a los datos de interlocking (personas)



31

Un grafo aleatorio ajustado a los datos de interlocking (personas)

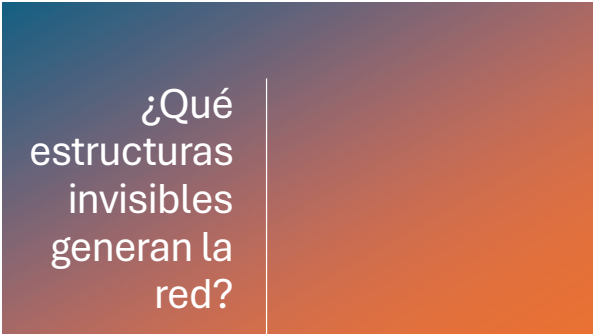


32

Referencias

- <https://statnet.org/trac>
- Buscar también p* models

33

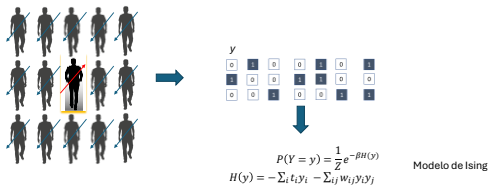


34

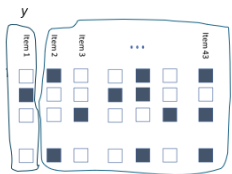


Ising model

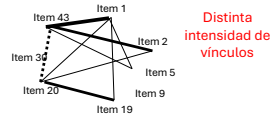
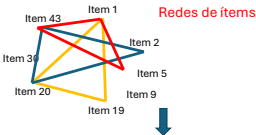
35



36



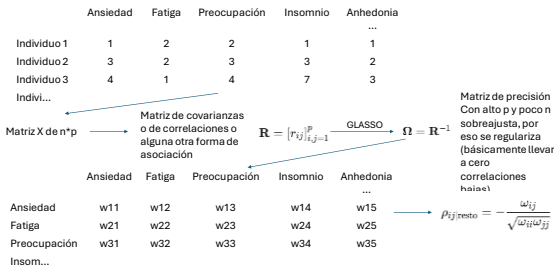
Nodos: Items son variables aleatorias
Links: Asociación condicional entre Items



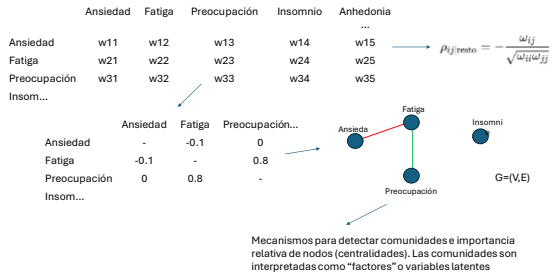
Distinta
intensidad de
vínculos

37

Exploratory Graph Analysis (EGA)



38



39

Ejemplo

$$p_{ij\text{rento}} = -\frac{\omega_{ij}}{\sqrt{\omega_i\omega_j}}$$

	Ashedonia	Incomiso	Fatiga	Irritabilidad	Ansiiedad	Preocupacion	Nerviosismo	Inquietud
Ashedonia	0.0000000	0.33826200	0.60978945	0.00000000	-0.01153948	-0.03252217	-0.01261822	-0.05895483
Incomiso	0.33826200	0.0000000	0.23114811	0.04261495	0.20439652	0.16532328	0.25728866	0.07648377
Fatiga	0.60978945	0.23114811	0.0000000	-0.06833799	0.00000000	0.00000000	0.00000000	-0.05248811
Irritabilidad	0.0000000	0.04261495	-0.06833799	0.00000000	0.21940004	0.17640801	0.23462142	0.22784654
Ansiiedad	-0.01153948	0.04261495	0.00000000	0.21940004	0.00000000	0.33021116	0.22713263	0.10862141
Preocupacion	-0.03252217	0.16532328	0.00000000	0.17640804	0.33021116	0.00000000	0.26209430	0.13346342
Nerviosismo	-0.01261822	0.25728866	0.00000000	0.23462143	0.22713263	0.26209430	0.00000000	0.28238449
Inquietud	-0.05895483	0.07648377	-0.05248809	0.22784654	0.10862141	0.13346342	0.28238449	0.00000000

40

Ejemplo

```
Model: G-LSD (LSD with gamma = 0.5)
Correlation data
Lambda: 0.000200120041182 (n = 100, ratio = 0.1)

Number of nodes: 8
Number of edges: 28
Edge density: 0.357

Non-zero edge weights:
  n   30   Min   Max
# 11 0.182 -0.889 0.618

Algorithm: Multitap
Number of communities: 2

Ashedonia  Incomiso  Fatiga  Irritabilidad  Ansiiedad  Preocupacion  Nerviosismo  Inquietud
1          1          1          1          1          1          1          1
2          2          2          2          2          2          2          2

Multidimensional Method: Linear
Multidimensional: No

LRF1: -0.576
```

41

Ejemplo

Configuración del modelo:
Indica la regularización hecha
Gamma – complejidad aceptada
Lambda – penalización utilizada

```
Model: G-LSD (LSD with gamma = 0.5)
Correlation data
Lambda: 0.000200120041182 (n = 100, ratio = 0.1)

Number of nodes: 8
Number of edges: 28
Edge density: 0.357

Non-zero edge weights:
  n   30   Min   Max
# 11 0.182 -0.889 0.618

Algorithm: Multitap
Number of communities: 2

Ashedonia  Incomiso  Fatiga  Irritabilidad  Ansiiedad  Preocupacion  Nerviosismo  Inquietud
1          1          1          1          1          1          1          1
2          2          2          2          2          2          2          2

Multidimensional Method: Linear
Multidimensional: No

LRF1: -0.576
```

42

Ejemplo

Descripción de la red

```
Model: G.6510 (CICIC with gamma = 0.5)
Correlations: none
Lambda: 0.000209120041182 (n = 100, ratio = 0.1)

Number of nodes: 8
Number of edges: 26
Edge density: 0.807

Non-zero edge weights:
  0  10  Min  Max
0.101 0.102 -0.009 0.018

-----
Algorithm: Multinet
Number of communities: 2
  Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1           1           1           1           1           1           1
2           2           2           2           2           2           2
Impurity: 2

-----
Classification Method: Louvain
Misclassification: 0%

-----
F1F2: -0.004
```

43

Ejemplo

Peso de las aristas

```
Model: G.6510 (CICIC with gamma = 0.5)
Correlations: none
Lambda: 0.000209120041182 (n = 100, ratio = 0.1)

Number of nodes: 8
Number of edges: 26
Edge density: 0.807

Non-zero edge weights:
  0  10  Min  Max
0.101 0.102 -0.009 0.018

-----
Algorithm: Multinet
Number of communities: 2
  Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1           1           1           1           1           1           1
2           2           2           2           2           2           2
Impurity: 2

-----
Classification Method: Louvain
Misclassification: 0%

-----
F1F2: -0.004
```

44

Ejemplo

Detección de comunidades y asignación de síntomas a las comunidad

```
Model: G.6510 (CICIC with gamma = 0.5)
Correlations: none
Lambda: 0.000209120041182 (n = 100, ratio = 0.1)

Number of nodes: 8
Number of edges: 26
Edge density: 0.807

Non-zero edge weights:
  0  10  Min  Max
0.101 0.102 -0.009 0.018

-----
Algorithm: Multinet
Number of communities: 2
  Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1           1           1           1           1           1           1
2           2           2           2           2           2           2
Impurity: 2

-----
Classification Method: Louvain
Misclassification: 0%

-----
F1F2: -0.004
```

45

Ejemplo

```
Model: 0.650 (CIC with gamma = 0.5)
Correlations: none
Lambda: 0.00020120041182 (n = 100, ratio = 0.1)
Number of nodes: 8
Number of edges: 26
Edge density: 0.807
Non-zero edge weights:
  0  30  Max
  0.101 0.102 -0.009 0.019
-----
Algorithm: Walktrap
Number of communities: 2
Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1          1          1          1          1          1          1
2          1          1          1          1          1          1
Depresion  2
-----
Modularity Q: 0.875
Stability S: 0.875
NMI: 0.875
```

Test de unidimensionalidad

46

Ejemplo

```
Model: 0.650 (CIC with gamma = 0.5)
Correlations: none
Lambda: 0.00020120041182 (n = 100, ratio = 0.1)
Number of nodes: 8
Number of edges: 26
Edge density: 0.807
Non-zero edge weights:
  0  30  Max
  0.101 0.102 -0.009 0.019
-----
Algorithm: Walktrap
Number of communities: 2
Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1          1          1          1          1          1          1
2          1          1          1          1          1          1
Depresion  2
-----
Modularity Q: 0.875
Stability S: 0.875
NMI: 0.875
```

Total EntropyFit Index: mientras más negativo mejor ajuste

47

Volviendo al punto del algoritmo walktrap...

- El 4to paso era detectar comunidades y nodos relevantes...

Detección de comunidades y asignación de síntomas a las comunidades

```
Model: 0.650 (CIC with gamma = 0.5)
Correlations: none
Lambda: 0.00020120041182 (n = 100, ratio = 0.1)
Number of nodes: 8
Number of edges: 26
Edge density: 0.807
Non-zero edge weights:
  0  30  Max
  0.101 0.102 -0.009 0.019
-----
Algorithm: Walktrap
Number of communities: 2
Anhedonia  Insomnia  Fatiga Irritabilidad  Ansiedad  Preocupacion  Nerviosismo
1          1          1          1          1          1          1
2          1          1          1          1          1          1
Depresion  2
-----
Modularity Q: 0.875
Stability S: 0.875
NMI: 0.875
```

48



Detección de comunidades

- Básicamente encontrar dónde están las densidades locales más fuertes y un criterio para separar nodos en subgrafos.
- El principal: Modularidad → maximizar Q:

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

1 si ambos están en la misma comunidad
0 en caso contrario

- Algoritmo más popular: Louvain

49



Detección de comunidades

- Otros algoritmos frecuentes
- Walktrap
 - Intuición: Caminatas aleatorias en el grafo deberían rescatar cuando te estás dando vuelta en un mismo vecindario.
 - Se calculan las distancias entre todos los pares de nodos usando caminatas aleatorias.
 - Los nodos cercanos según esta distancia se agrupan iterativamente, creando una jerarquía de posibles particiones.
- Infomap
 - Objetivo: encontrar la partición de la red que comprime mejor la información de la caminata aleatoria

50
