

MODULARITY AND EDGE SAMPLING

Fiona Skerman
Bristol/Uppsala

Birmingham 2020

PARTITIONING NETWORKS:

Network:

links between blogs on climate change

ELGESEM, STESKAL & DIAKOPOULOS 2015

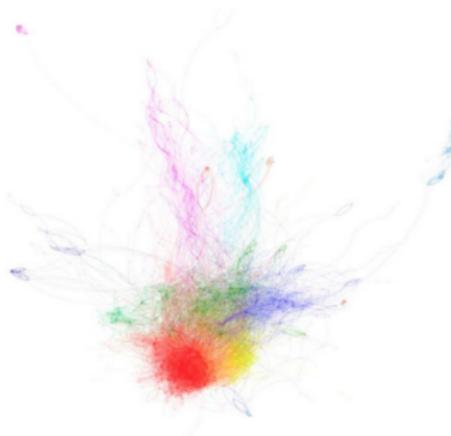


Figure 1. The network of climate change blogs, colored by community.

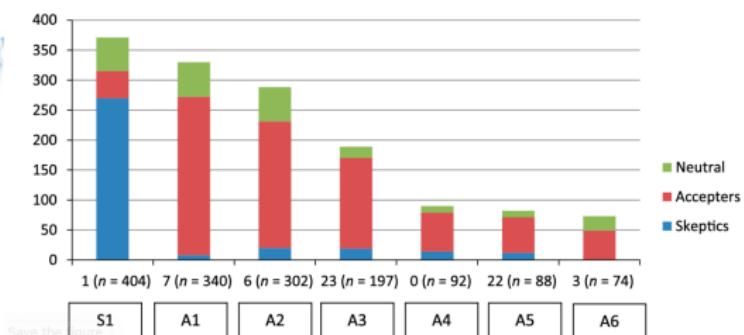


Figure 3. The distribution of skeptical, accepting, and neutral blogs in the seven largest among the central groups of blogs concerned with climate change.

PARTITIONING NETWORKS:

Network:

links between blogs on climate change

ELGESEM, STESKAL & DIAKOPOULOS 2015

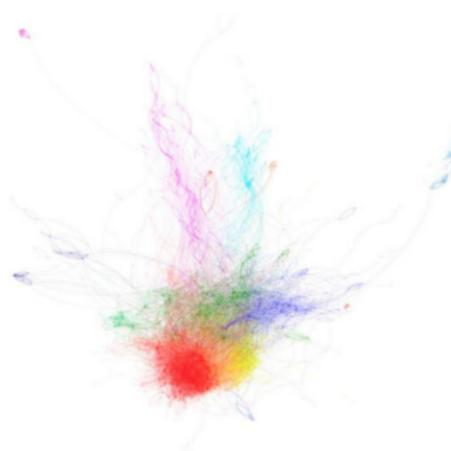


Figure 1. The network of climate change blogs, colored by community.

Table 5. The top 15 collocates around "climate" in communities 1 (skeptic), 23 (accepter), and 7 (accepter) computed with the point-wise mutual information metric.

Top collocates of "CLIMATE" in the skeptical community S1	Top collocates of "CLIMATE" in the accepter community A3	Top collocates of "CLIMATE" in the accepter community A1
1 CLIMATE	1 DENIERS	1 POPPIN
2 SKEPTICS	2 SKEPTICS	2 DENIERS
3 ALARMISM	3 CLIMAT	3 SKEPTICS
4 DENIERS	4 DECADAL	4 OBAMA
5 IPCC	5 CONTRARIANS	5 WWW
6 DECADAL	6 OBAMA	6 EU'S
7 ALARMISTS	7 NOAA'S	7 CLIMATE
8 CLIMAT	8 AGW	8 YVO
9 CHANGE	9 WWW	9 NOAA'S
10 INTERGOVERNMENTAL	10 DENIER	10 WILDFIRES
11 OBAMA	11 CLIMATE	11 CHANGE'S
12 ANTHROPOGENIC	12 VAPOR	12 IPCC
13 AGW	13 ANTHROPOGENIC	13 ALARMISM
14 IPCC'S	14 ALARMISM	14 PACHAURI
15 WARMING	15 CONTRARIAN	15 DENIER

Reference corpus: The British National Corpus, approximately 100 million words.

DEFINITION

Motivation - configuration model

NEWMAN & GIRVAN 2004

Let G be a graph on $m \geq 1$ edges and \mathcal{A} a vertex partition,

Write d_u for the degree of u ,

$e(A)$ for the number of edges in A and $\text{vol}(A) = \sum_{u \in A} d_u$.

$$q_{\mathcal{A}}(G) = \frac{1}{2m} \sum_{A \in \mathcal{A}} \sum_{u, v \in A} \left(1_{uv \in E} - \frac{d_u d_v}{2m}\right) = \sum_{A \in \mathcal{A}} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{4m^2}$$



Expected number of edges between distinct u and v is $\frac{d_u d_v}{2m - 1}$

DEFINITION

Motivation - configuration model

NEWMAN & GIRVAN 2004

Let G be a graph on $m \geq 1$ edges and \mathcal{A} a vertex partition,

Write d_u for the degree of u ,

$e(A)$ for the number of edges in A and $\text{vol}(A) = \sum_{u \in A} d_u$.

$$q_{\mathcal{A}}(G) = \frac{1}{2m} \sum_{A \in \mathcal{A}} \sum_{u, v \in A} \left(1_{uv \in E} - \frac{d_u d_v}{2m} \right) = \sum_{A \in \mathcal{A}} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{4m^2}$$

Max. Modularity

$$q^*(G) = \max_{\mathcal{A}} q_{\mathcal{A}}(G)$$

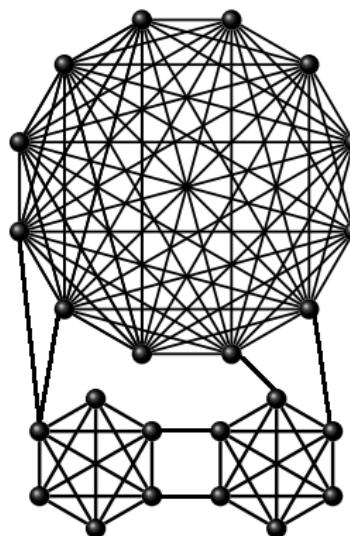
Edge contribution

$$q_{\mathcal{A}}^E(G) = \sum_{A \in \mathcal{A}} \frac{e(A)}{m}$$

Degree tax

$$q_{\mathcal{A}}^D(G) = \sum_{A \in \mathcal{A}} \frac{\text{vol}(A)^2}{4m^2}$$

Example Graph



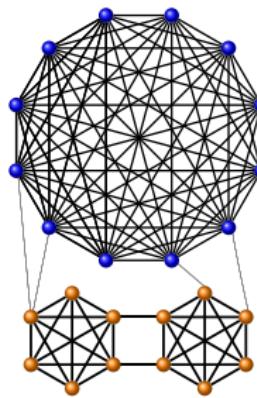
Edge contribution

$$q_{\mathcal{A}}^E(G) = \sum_{A \in \mathcal{A}} \frac{e(A)}{m}$$

Degree tax

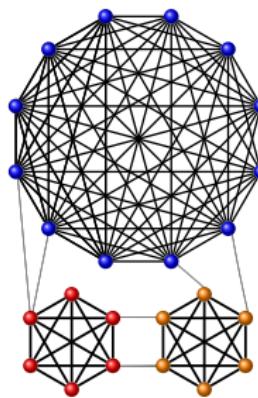
$$q_{\mathcal{A}}^D(G) = \sum_{A \in \mathcal{A}} \frac{\text{vol}(A)^2}{4m^2}$$

3 Possible Partitions



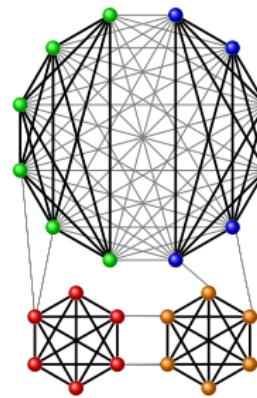
$$q_{\mathcal{A}_1}^E = 0.96, \quad q_{\mathcal{A}_1}^D = 0.56$$

$$q_{\mathcal{A}_1} = 0.40$$



$$q_{\mathcal{A}_2}^E = 0.94, \quad q_{\mathcal{A}_2}^D = 0.50$$

$$q_{\mathcal{A}_2} = 0.44$$



$$q_{\mathcal{A}_3}^E = 0.59, \quad q_{\mathcal{A}_3}^D = 0.29$$

$$q_{\mathcal{A}_3} = 0.30$$

PROPERTIES OF MODULARITY

Edge contribution

$$q_{\mathcal{A}}^E(G) = \sum_{A \in \mathcal{A}} \frac{e(A)}{m}$$

Degree tax

$$q_{\mathcal{A}}^D(G) = \sum_{A \in \mathcal{A}} \frac{\text{vol}(A)^2}{4m^2}$$

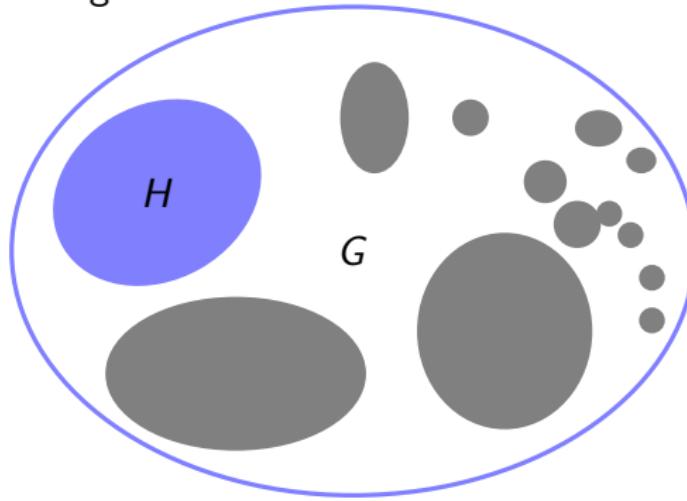
Basic Properties

- $0 \leq q^*(G) < 1.$
- ‘blind’ to isolated vertices. (V_0)
- parts connected in an optimal partition.
 $A \setminus V_0$ connected for all $A \in \mathcal{A}$ if $q_{\mathcal{A}}(G) = q^*(G)$.
- Resolution limit, Robustness.

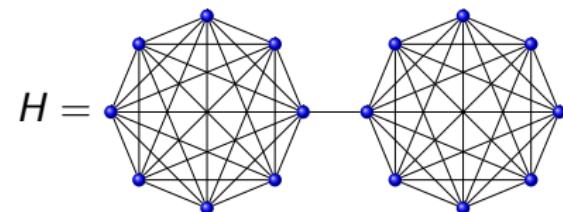
RESOLUTION LIMIT

FORTUNATO AND BARTHÉLEMY 08

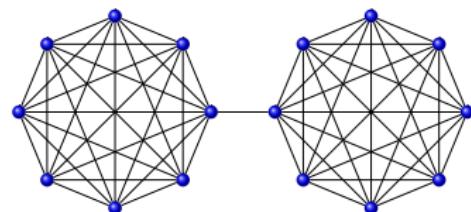
Subgraph H
 h edges



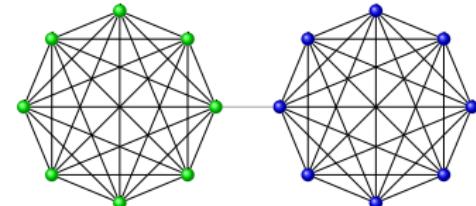
Graph G , m edges



$$H =$$



If $h < \sqrt{2m}$, e.g. $m = 1625$.



If $h > \sqrt{2m}$, e.g. $m = 1624$.

ROBUSTNESS

Partition structure is sensitive to noise in edges

Modularity value is robust

PROPOSITION (McDIARMID, S.)

Let $G = (V, E)$ and $G' = (V, E')$ with $|E| \geq |E'|$,

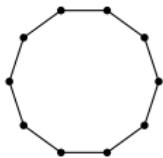
$$|q^*(G) - q^*(G')| \leq \frac{2|E \setminus E'|}{|E|}$$

High k cliques

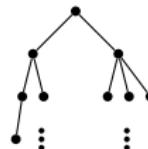
$$q^* \rightarrow 1$$



$$q^* = 1 - \frac{1}{k}$$

cycle C_n 

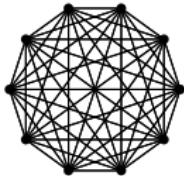
$$q^* \sim 1 - \frac{2}{\sqrt{n}}$$

tree, max deg Δ 

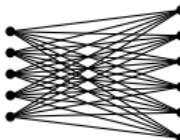
$$q^* \geq 1 - 4\sqrt{\frac{\Delta}{n}}$$

Lowclique or near-clique
 $K_n \setminus E'$, $|E'| \leq \frac{n}{2}$

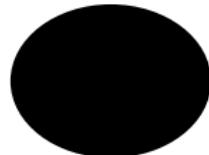
$$q^* \rightarrow 0$$



$$q^* = 0$$

complete
multipartite

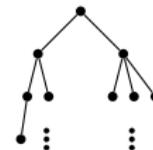
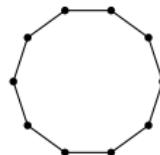
$$q^* = 0$$

 $G(n, \frac{\log n}{n})$ random graph
connectivity threshold

$$q^* = O\left(\frac{1}{\sqrt{\log n}}\right)$$

High

$$q^* \rightarrow 1$$

 **k cliques****cycle C_n** **tree, max deg Δ** 

$$q^* = 1 - \frac{1}{k}$$

$$q^* \sim 1 - \frac{2}{\sqrt{n}}$$

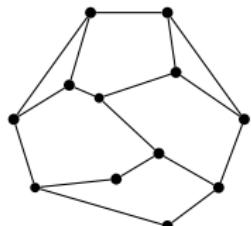
$$q^* \geq 1 - 4\sqrt{\frac{\Delta}{n}}$$

Critical

random cubic

random r -regular
large r $G(n, \frac{c}{n})$ random graph
large $c > 1$

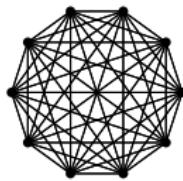
$$\epsilon < q^* < 1 - \epsilon$$



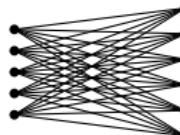
$$0.66 < q^* < 0.81$$

**Low**clique or near-clique
 $K_n \setminus E'$, $|E'| \leq \frac{n}{2}$ complete
multipartite $G(n, \frac{\log n}{n})$ random graph
connectivity threshold

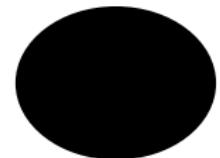
$$q^* \rightarrow 0$$



$$q^* = 0$$



$$q^* = 0$$



$$q^* = O\left(\frac{1}{\sqrt{\log n}}\right)$$

SAMPLING TO DETERMINE MODULARITY

Let G be the underlying network

Let G_p be obtained by sampling each edge with probability p .

How does $q^*(G_p)$ behave? For which p do we have $q^*(G_p) \approx q^*(G)$?

PHASE TRANSITION IN ERDŐS-RÉNYI $(K_n)_p$

THEOREM (McDIARMID, S.)

- (i) If $n^2 p \rightarrow \infty$ and $np \leq 1 + o(1)$ then whp $q^*(G_n) \rightarrow 1$. SPARSE
- (ii) If $np \rightarrow c > 1$ then $\exists \epsilon = \epsilon(c) > 0$ such that whp
 $\epsilon < q^*(G_n) < 1 - \epsilon$. CRITICAL
- (iii) If $np \rightarrow \infty$ then whp $q^*(G_n) \rightarrow 0$. DENSE

THEOREM (MC DIARMID, S.)

There exists $0 < a < b$ such that, if $1/n \leq p = p(n) \leq 0.99$ then

$$\frac{a}{\sqrt{np}} < q^*(G_n) < \frac{b}{\sqrt{np}} \quad \text{whp.}$$

SAMPLING TO DETERMINE MODULARITY

Let G be the underlying network

Let G_p be obtained by sampling each edge with probability p .

THEOREM (McDIARMID, S.)

$\forall \varepsilon > 0, \exists c$ such that the following holds with probability at least $1 - \varepsilon$ if $0 < p < 1$ and graph G satisfy $e(G)p \geq c$, then

$$q^*(G_p) > q^*(G) - \varepsilon.$$

if $0 < p < 1$ and graph G satisfy $e(G)p \geq cn$, then

$$q^*(G_p) < q^*(G) + \varepsilon.$$

SAMPLING TO DETERMINE MODULARITY

Let G be the underlying network

Let G_p be obtained by sampling each edge with probability p .

THEOREM (McDIARMID, S.)

$\forall \varepsilon > 0, \exists c$ such that the following holds with probability at least $1 - \varepsilon$ if $0 < p < 1$ and graph G satisfy $e(G)p \geq c$, then

$$q^*(G_p) > q^*(G) - \varepsilon.$$

if $0 < p < 1$ and graph G satisfy $e(G)p \geq cn$, then

$$q^*(G_p) < q^*(G) + \varepsilon.$$

Observe: if $e(G)p \leq c_0$ then $\mathbb{P}(e(G_p) = 1) \geq c_0 e^{-2c_0}$.

Earlier results: $q^*(K_n) = 0$, for $e(K_n)p \leq c_0$ whp $q^*((K_n)_p) \geq \epsilon(c_o)$.

IDEA OF PROOF

Say vertex partition \mathcal{A} is η -fat for G if each part has volume $\geq \eta \text{ vol}(G)$.

For non-empty G and vertex partition \mathcal{A}_0 , can construct η -fat 'coarser' \mathcal{A} with

$$q_{\mathcal{A}}(G) > q_{\mathcal{A}_0}(G) - 2\eta.$$

I): If $0 < p < 1$ and $e(G)p \geq c$, then w.p. $1 - \epsilon$, $q^*(G_p) > q^*(G) - \epsilon$.

Can assume $q^*(G) = q_{\mathcal{A}_0}(G) \geq \epsilon$.

From \mathcal{A}_0 construct 'coarser' partition \mathcal{A} $\frac{\epsilon}{4}$ -fat for G .

Show $q_{\mathcal{A}}(G_p) > q_{\mathcal{A}}(G) - \frac{\epsilon}{2}$.

Edges internal to \mathcal{A} in G is at least $\frac{\epsilon}{2}e(G)$.

Edges internal to \mathcal{A} in G_p approx $\frac{\epsilon}{2}e(G)p$. Thus $q_{\mathcal{A}}^E(G_p) > q_{\mathcal{A}}^E(G) - \frac{\epsilon}{4}$.

IDEA OF PROOF

Say vertex partition \mathcal{A} is η -fat for G if each part has volume $\geq \eta \text{ vol}(G)$.

For non-empty G and vertex partition \mathcal{A}_0 , can construct η -fat 'coarser' \mathcal{A} with

$$q_{\mathcal{A}}(G) > q_{\mathcal{A}_0}(G) - 2\eta.$$

II): If $0 < p < 1$ and $e(G)p \geq cn$, then w.p. $1 - \epsilon$, $q^*(G_p) < q^*(G) + \epsilon$.

Define bad events:

$\mathcal{B}_1 \exists A \subseteq V$ with $\text{vol}(G) < \frac{\eta}{2} \text{vol}(G)$ and $\text{vol}(G_p) > \eta \text{vol}(G_p)$.

i.e. $\exists A$ not $\frac{\eta}{2}$ -fat for G but is η -fat for G_p .

$\mathcal{B}_2 \exists \mathcal{A}$ $\frac{\eta}{2}$ -fat for G and $q_{\mathcal{A}}(G_p) < q_{\mathcal{A}}(G) + \epsilon/2$.

EDGE LIMITED SEARCH

c-edge-limited search operates as follows.

initialise the sets F to be empty and \tilde{E} to contain all the edges of K_n .

while $\tilde{E} \neq \emptyset$ and $|F| < c$,

sample $e \in \tilde{E}$ uniformly at random, and delete it from \tilde{E}

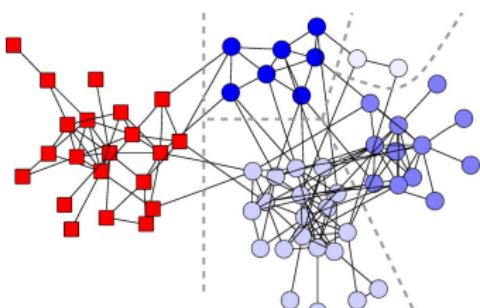
if $e \in E(G)$ **then** add e to F

return the graph $G^c = (V, F)$, where $V = [n]$.

COROLLARY (McDIARMID, S.)

Given $\epsilon > 0$ there exists $c = c(\epsilon)$ such that the following holds. When we run edge-limited search on G , with probability $> 1 - \epsilon$,
the random output graph G^c satisfies $q^*(G^c) > q^*(G) - \epsilon$
the random output graph G^{cn} satisfies $q^*(G^{cn}) < q^*(G) + \epsilon$

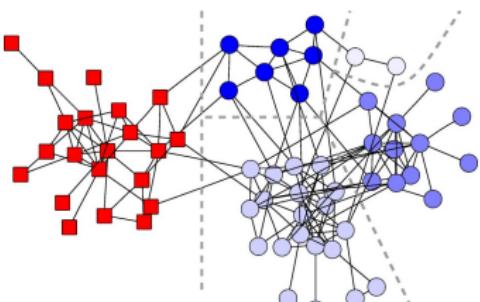
LUSSEAU PHD THESIS



dolphins = 62
edges = 159

$q^* = 0.52$

LUSSEAU PHD THESIS



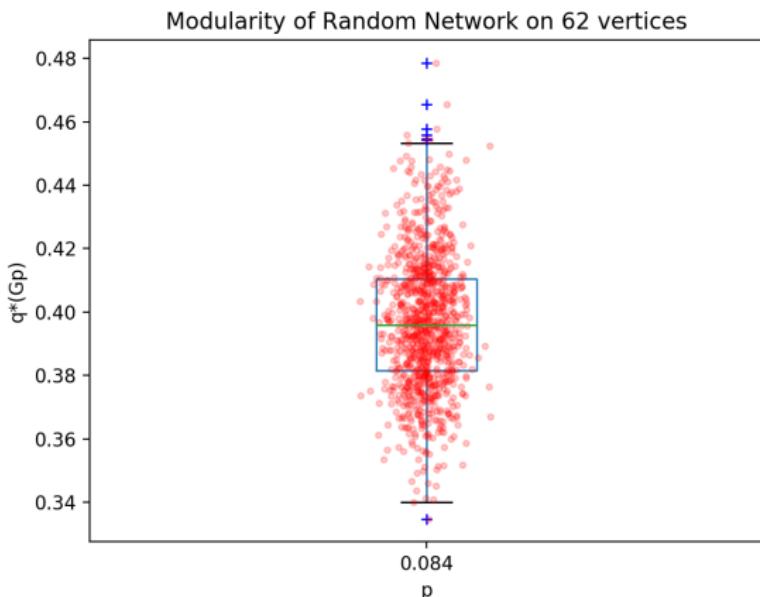
dolphins = 62
edges = 159

$q^* = 0.52$

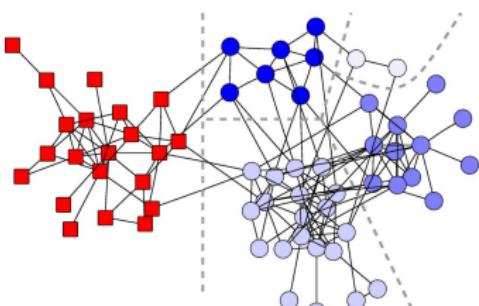
8.4% of possible edges

RANDOM DATA

$$q^*(\text{dolphins}) > q^*(\text{random network})???$$



LUSSEAU PHD THESIS



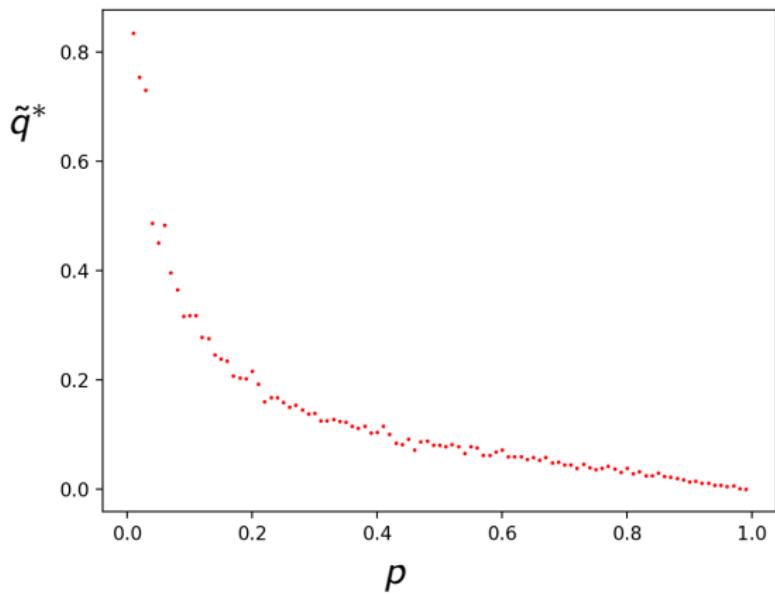
dolphins = 62
edges = 159

$q^* = 0.52$

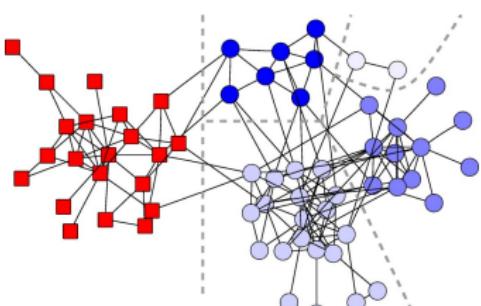
8.4% of possible edges

RANDOM DATA

Simulate 62 vertices, with edge prob p .



LUSSEAU PHD THESIS



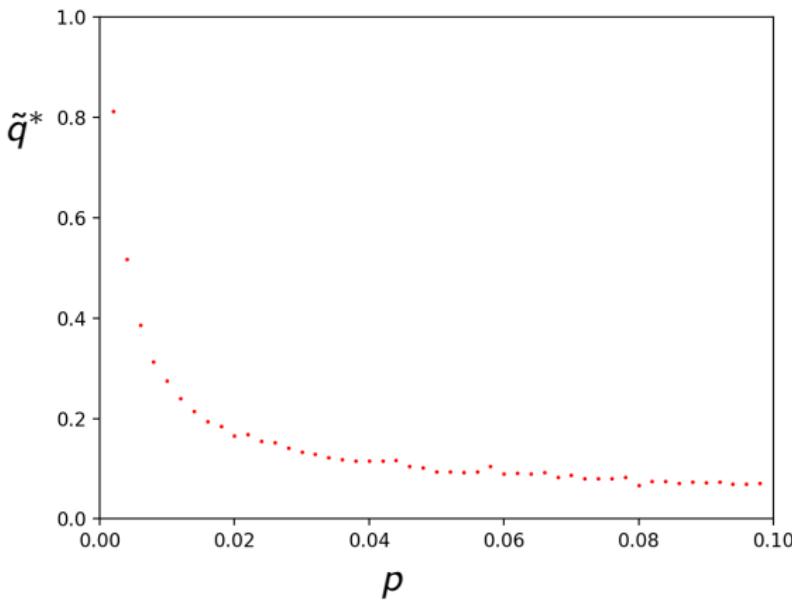
dolphins = 62
edges = 159

$$q^* = 0.52$$

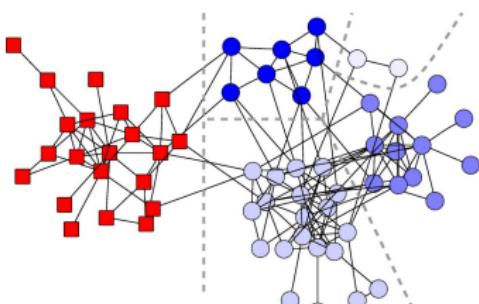
8.4% of possible edges

RANDOM DATA

Simulate 1000 vertices, with edge prob p .



LUSSEAU PHD THESIS



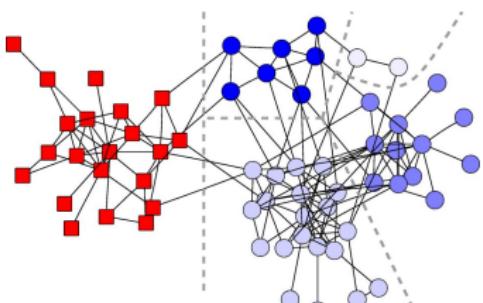
dolphins = 62
edges = 159

$q^* = 0.52$

UNDERLYING NETWORK??

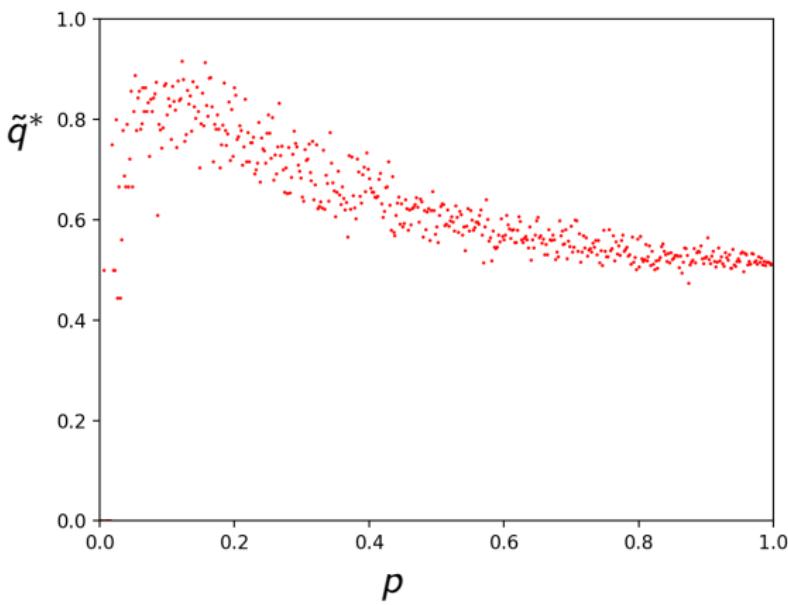
$q^*(\text{observed})$ vs $q^*(\text{underlying})$?

LUSSEAU PHD THESIS

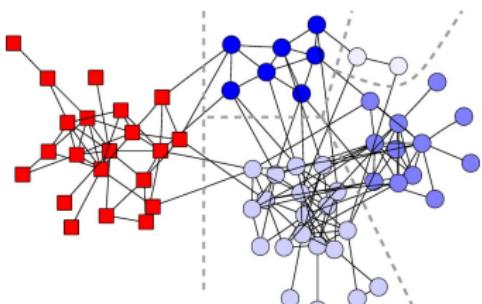


dolphins = 62
edges = 159

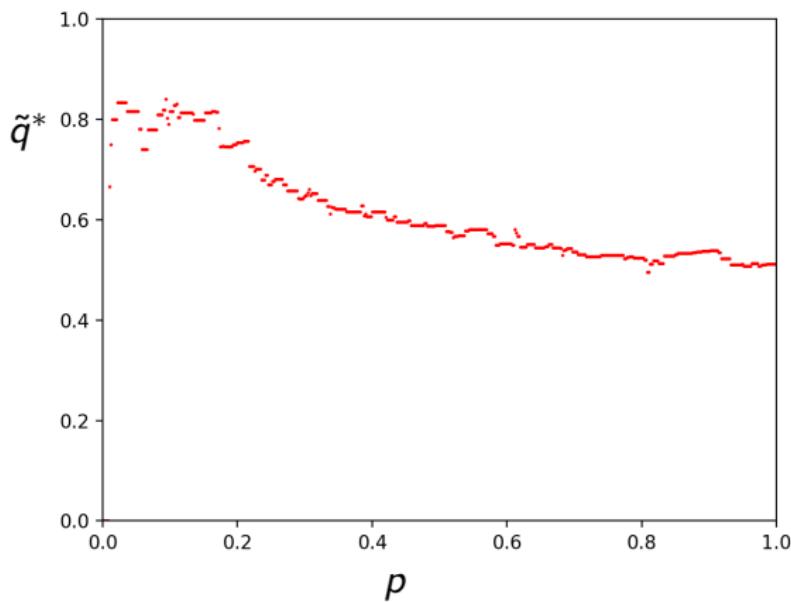
$$q^* = 0.52$$

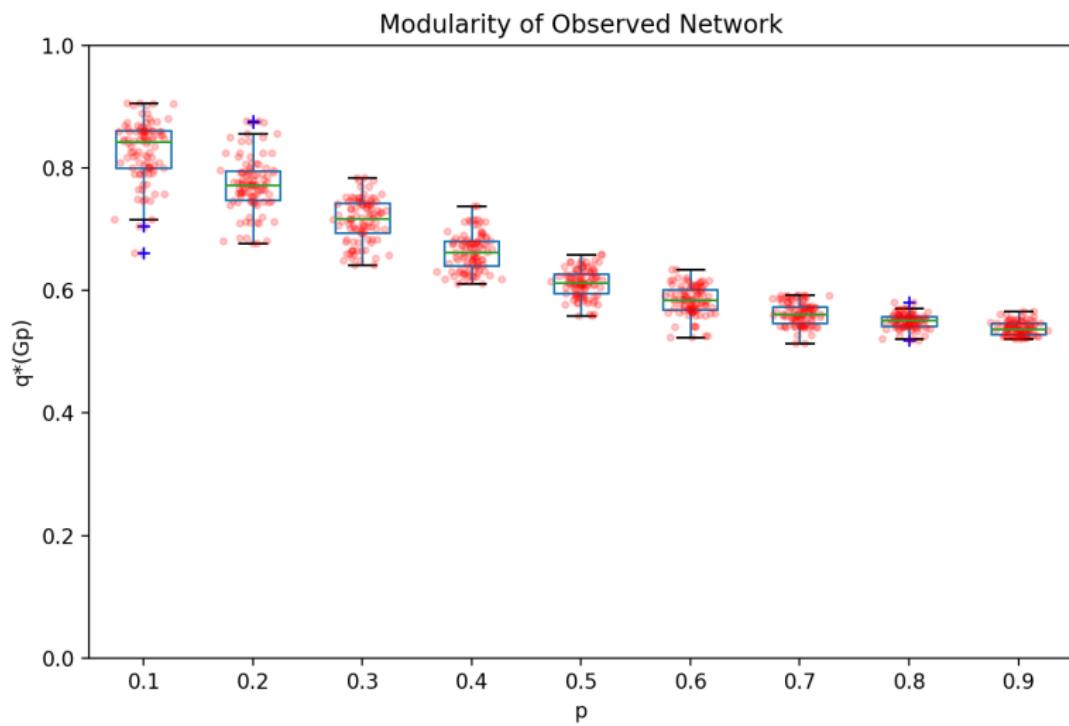


LUSSEAU PHD THESIS



$$q^* = 0.52$$





Ad: PARAMETRISED COMPLEXITY OF MODULARITY

THEOREM (MEEKS, S.)

Computing the modularity is w[1]-hard when parametrised by pathwidth.

OPEN QUESTION:

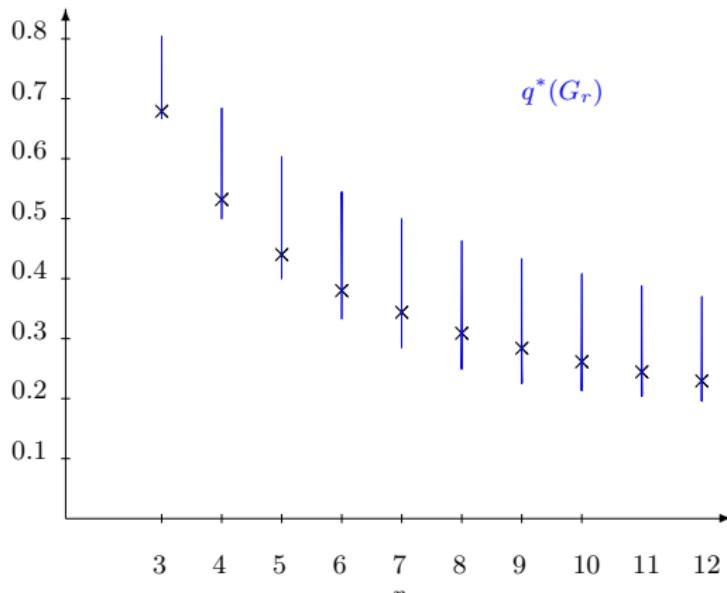
Parametrised by shrub-depth?

Idea for hardness result

Suppose only allowed partitions where each part has exactly 4 crossing-edges. Then over any possible graphs and partitions, get **best modularity** if $\text{vol}(A) = 2\sqrt{2m}$ for each part. Reduces to AECP [ENCISO ET AL. 2009]

RANDOM r -REGULAR GRAPHS

	$r =$	3	4	5	6	7	8	9	10	11	12
bounds	$q^*(G_r) >$	0.666	0.500	0.400	0.333	0.285	0.250	0.226	0.214	0.204	0.196
	$q^*(G_r) <$	0.804	0.684	0.603	0.544	0.499	0.463	0.433	0.408	0.388	0.370
simulations	louvain	0.679	0.531	0.440	0.380	0.343	0.312	0.284	0.262	0.244	0.230
	reshuffle	0.677	0.531	0.446	0.391	0.353	0.326	0.303	0.285	0.269	0.256



McDIARMID & S.
PROKHORENKOVA ET AL

RANDOM r -REGULAR GRAPHS

	$r =$	3	4	5	6	7	8	9	10	11	12
bounds	$q^*(G_r) >$	0.666	0.500	0.400	0.333	0.285	0.250	0.226	0.214	0.204	0.196
	$q^*(G_r) <$	0.804	0.684	0.603	0.544	0.499	0.463	0.433	0.408	0.388	0.370
simulations	louvain	0.679	0.531	0.440	0.380	0.343	0.312	0.284	0.262	0.244	0.230
	reshuffle	0.677	0.531	0.446	0.391	0.353	0.326	0.303	0.285	0.269	0.256

OPEN QUESTION:

Are random graphs whp the least modular r -regular graphs?

$$q_r^-(n) = \min_{G \in \mathcal{G}_{r,n}} q^*(G)$$

KNOWN:

$$\begin{aligned} q_2^-(n) &= 1 - 2.041/\sqrt{n} \\ q^*(G_2) &= 1 - 2/\sqrt{n} \text{ whp} \end{aligned}$$

For $r > 2$,

$$q_r^-(n) \geq 2/r - o(1)$$

OPEN:

Does there exist $\delta > 0$ such that
 $q^*(G_3) > 2/3 + \delta$ whp ?

For $r > 2$ is it true that

$$q_r^-(n) = q^*(G_r)(1 + o(1)) \text{ whp ?}$$

OPEN QUESTIONS

Consider $G_{n,c/n}$ for large $c > 1$.

Is it true whp $q^*(G_{n,c/n}) = \frac{0.97}{\sqrt{c}} + o_c(\frac{1}{\sqrt{c}})$?

Is 5 communities best whp? (no)

Does there exist a positive integer k (perhaps $k = 5$?) with the property that, for each $\epsilon > 0$ there exists c_0 such that, if $c \geq c_0$ then whp

$$q_{\leq k}(G_{n,c/n}) \geq q^*(G_{n,c/n})(1 - \epsilon).$$

