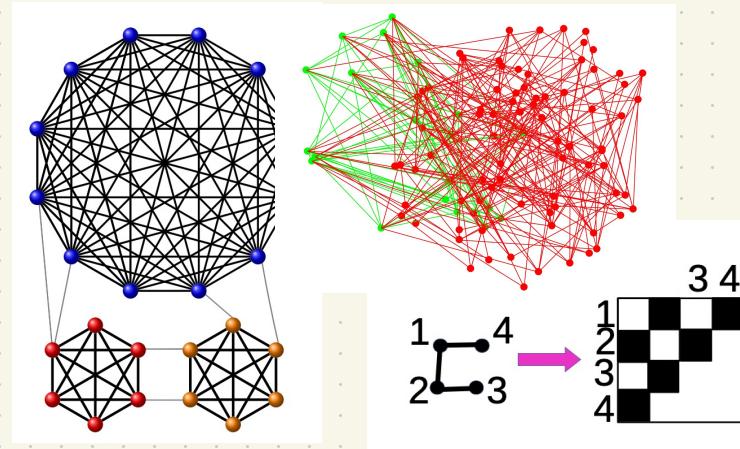
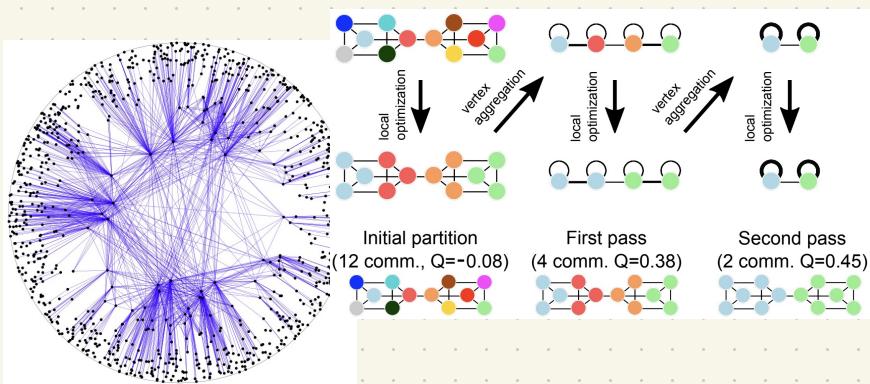


STRUCTURE IN NOISY NETWORKS

Fiona Skerman

- I: Modularity based clustering
- II: Fundamental limits of learning



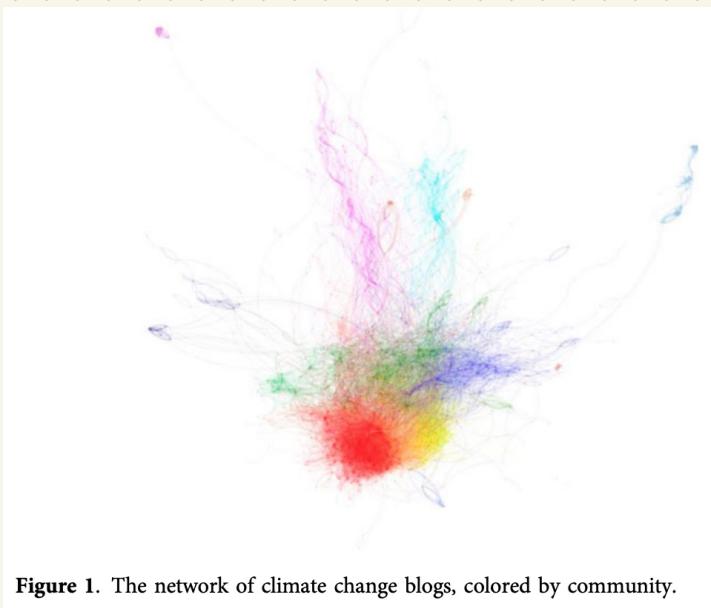
I: Modularity-based clustering

- Louvain: most commonly used method of community detection
- modularity: definitions + properties.

Example : Linguistics

V = 300 climate change blogs

E ~ based on links between blogs



Elgesam D , Steskal L. + Diakopoulos

"Structure and content of the discourse on climate change in the blogosphere"

Environmental Communication '2015 .

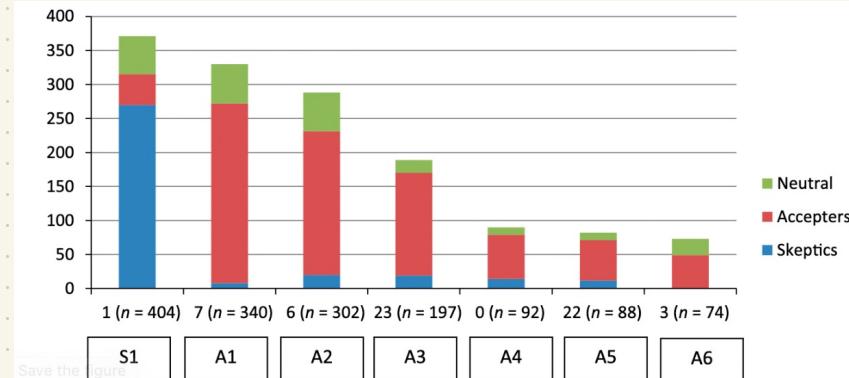


Figure 3. The distribution of skeptical, accepting, and neutral blogs in the seven largest among the central groups of blogs concerned with climate change.

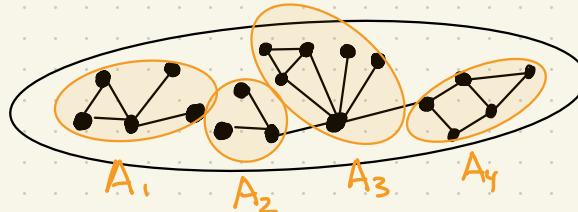
Table 5. The top 15 collocates around "climate" in communities 1 (skeptic), 23 (accepter), and 7 (accepter) computed with the point-wise mutual information metric.

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.

Modularity 'meas. of how well a graph can be clustered'

G



$$A = \{A_1, \dots, A_k\}$$

graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

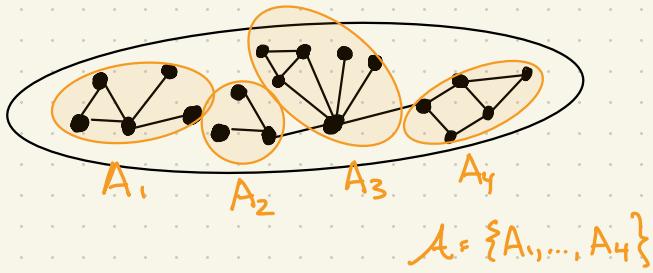
score of partition A , $q_A(G) =$

modularity of G $q^*(G) = \max_A q_A(G)$

'higher values taken to indicate
more community structure'

Modularity 'meas. of how well a graph can be clustered'

G



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G $q^*(G) = \max_A q_A(G)$

'higher values taken to indicate
more community structure'

Community Detection

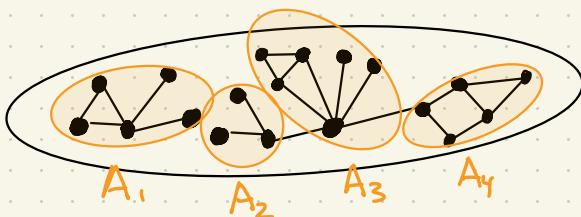
input graph $G = (V, E)$
 vertices nodes
 edges (weighted)

output vertex partition A
 'community division'

- modularity score NP-hard to opt.
- Louvain ~ modularity based
 & Leiden ~ iteratively build a partition
 local choices - maximise mod.
- most popular methods use modularity

Modularity 'meas of how well a graph can be clustered'

G



$$A = \{A_1, \dots, A_4\}$$

graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbb{1}_{[u \sim v]} - \frac{d_u \cdot d_v}{2m}$$

"edge contrib." "degree tax"

$$\rightarrow \frac{|A|^2}{n^2} \text{ for regular graphs}$$

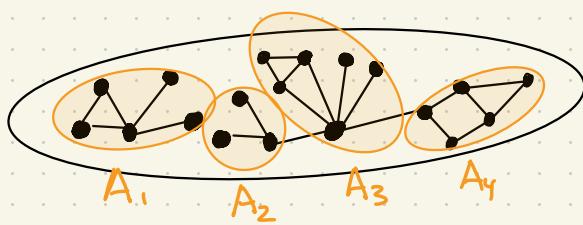
$d_u = \# \text{edges incident to } u$

$\text{vol}(A) = \# \text{edges in set } A$

$$\text{vol}(A) = \sum_{u \in A} d_u$$

Modularity 'meas of how well a graph can be clustered'

G



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbb{1}_{[u \sim v]} - \frac{d_u d_v}{2m}$$

"edge contrib." "degree tax"

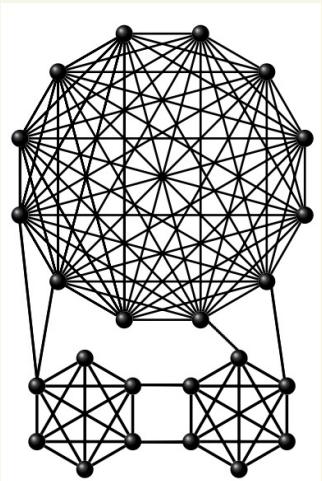
$d_u = \# \text{edges incident to } u$

$\text{vol}(A) = \# \text{edges in set } A$

$$\text{vol}(A) = \sum_{u \in A} d_u$$

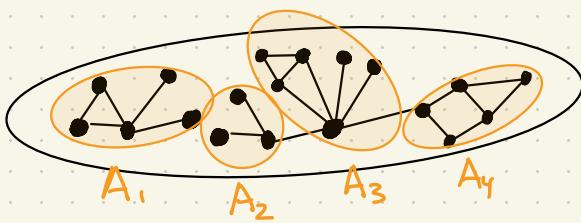
$$\rightarrow \frac{|A|^2}{n^2} \text{ for regular graphs}$$

Example



Modularity 'meas of how well a graph can be clustered'

G



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbf{1}_{[u \sim v]} - \frac{d_u \cdot d_v}{2m}$$

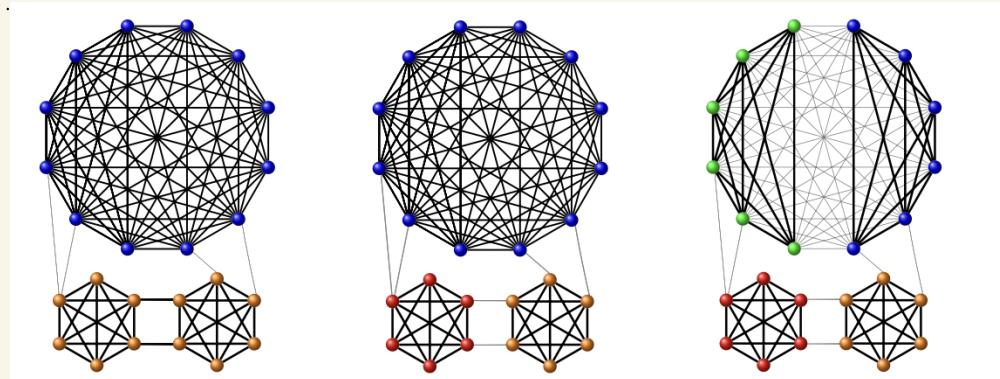
"edge contrib." "degree tax"

$d_u = \# \text{edges incident to } u$

$\text{vol}(A) = \# \text{edges in set } A$

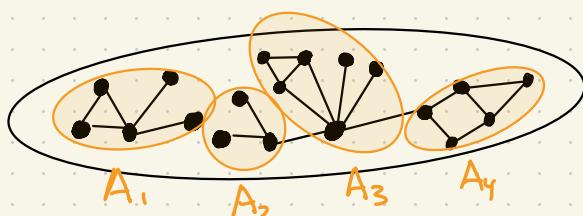
$$\text{vol}(A) = \sum_{u \in A} d_u$$

Example



Modularity 'meas of how well a graph can be clustered'

G



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbf{1}_{[u \sim v]} - \frac{d_u \cdot d_v}{2m}$$

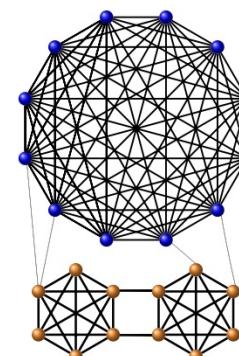
"edge contrib." "degree tax"

$d_u = \# \text{edges incident to } u$

$\text{vol}(A) = \# \text{edges in set } A$

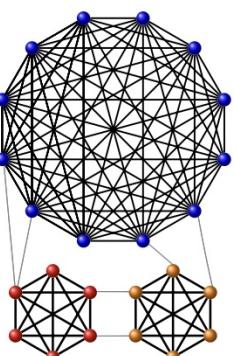
$$\text{vol}(A) = \sum_{u \in A} d_u$$

Example



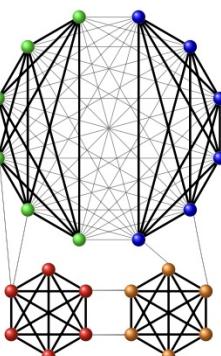
$$q_{A_1}^E = 0.96, \quad q_{A_1}^D = 0.56$$

$$q_{A_1} = 0.40$$



$$q_{A_2}^E = 0.94, \quad q_{A_2}^D = 0.50$$

$$q_{A_2} = 0.44$$

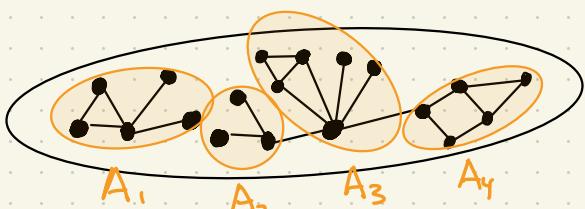


$$q_{A_3}^E = 0.59, \quad q_{A_3}^D = 0.29$$

$$q_{A_3} = 0.30$$

Modularity 'meas of how well a graph can be clustered'

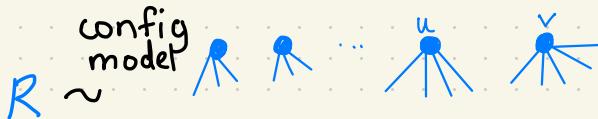
G



$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbb{1}_{[u \sim v]} - \frac{d_u \cdot d_v}{2m}$$

"edge contrib." "degree tax"



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

high vals taken to indicate
more community structure

$d_u = \# \text{edges incident to } u$

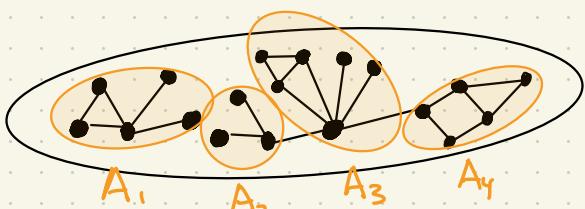
$\text{vol}(A) = \# \text{edges in set } A$

$$\text{vol}(A) = \sum_{u \in A} d_u$$

$$q_A(G) \approx \frac{1}{m} (e_G^{\text{int}}(A) - \mathbb{E}[e_R^{\text{int}}(A)])$$

Modularity 'meas of how well a graph can be clustered'

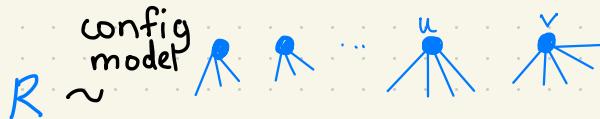
G



$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbb{1}_{[u \sim v]} - \frac{d_u d_v}{2m}$$

"edge contrib." "degree tax"



$$\cdot u \approx v \quad \mathbb{E}[\# \text{edges } u \sim v \text{ in } R] = \frac{d_u d_v}{2m-1}$$

graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

"high vals taken to indicate
more community structure"

$d_u = \# \text{edges incident to } u$

$\text{vol}(A) = \# \text{edges in set } A$

$$\text{vol}(A) = \sum_{u \in A} d_u$$

$$q_A(G) \approx \frac{1}{m} (e_G^{\text{int}}(A) - \mathbb{E}[e_R^{\text{int}}(A)])$$

$$\cdot \mathbb{E}[\# \text{edges within parts of } A \text{ in } R] = \sum_{A \in A} \frac{\text{vol}(A)(\text{vol}(A)-1)}{2m(2m-1)}$$

Fast unfolding of communities in large networks

[VD Blondel](#), [JL Guillaume](#), [R Lambiotte](#)... - *Journal of statistical ...*, 2008 - iopscience.iop.org

We propose a simple method to extract the community structure of large networks. Our method is a heuristic method that is based on modularity optimization. It is shown to outperform all ...

[☆ Save](#) [🔗 Cite](#) [Cited by 25562](#) [Related articles](#) [All 35 versions](#)

forest mice



About 410 results (0.07 sec)

Fast unfolding of communities in large networks

Search within citing articles

[HTML] Molecular logic of cellular diversification in the **mouse** cerebral cortex

[DJ Di Bella](#), [E Habibi](#), [RR Stickels](#), [G Scalia](#), [J Brown](#)... - *Nature*, 2021 - nature.com

... To define the closest identity of these cells, we applied a multi-class random **forest** classifier trained on the wild-type cell types (Extended Data Fig. 11f). Most of the knockout cells were ...

[☆ Save](#) [🔗 Cite](#) [Cited by 360](#) [Related articles](#) [All 11 versions](#)

[PDF] Identification of protein functions in **mouse** with a label space partition method

[X Li](#), [L Lu](#), [L Chen](#) - *Math. Biosci. Eng.*, 2022 - aimspress.com

... This study presented a new multi-label classifier for identifying functions of **mouse** proteins. Due to the number of functional types, which were termed as labels in the classification ...

[☆ Save](#) [🔗 Cite](#) [Cited by 31](#) [Related articles](#) [All 5 versions](#)

[HTML] Molecular architecture of the **mouse** nervous system

[A Zeisel](#), [H Hochgerner](#), [P Lönnerberg](#), [A Johnsson](#)... - *Cell*, 2018 - cell.com

... To assess the robustness of the clusters, we trained a random **forest** classifier to recognize cluster labels and then assessed its performance on held-out data (80% training set, 20% test ...

[☆ Save](#) [🔗 Cite](#) [Cited by 2519](#) [Related articles](#) [All 22 versions](#)

[HTML] Decoding molecular and cellular heterogeneity of **mouse** nucleus accumbens

[R Chen](#), [TR Blosser](#), [MN Djekidel](#), [J Hao](#)... - *Nature* ..., 2021 - nature.com

... In this study, we generated a cell census of the **mouse** NAc using single-cell RNA sequencing and multiplexed error-robust fluorescence *in situ* hybridization, revealing a high level of ...

RESEARCH PAPER

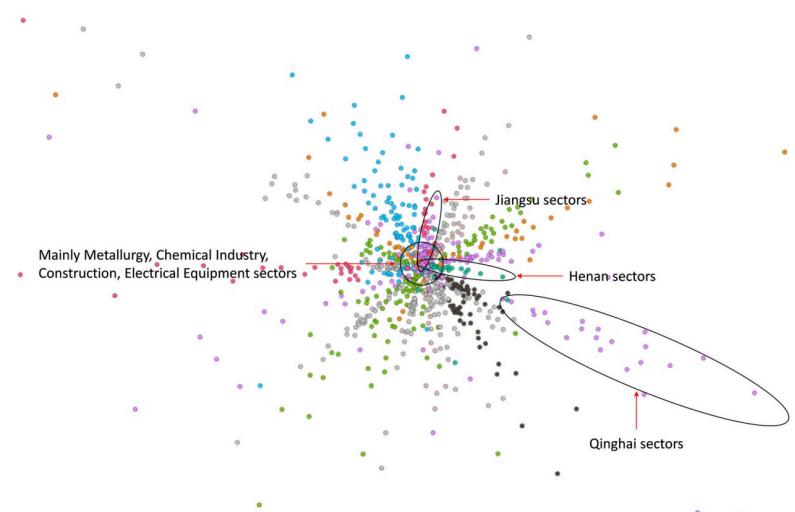
Hub disruption in patients with chronic neck pain: a graph analytical approach

[De Pauw](#), [Robby](#)^{a,*}; [Aerts](#), [Hannelore](#)^b; [Siugzdaitė](#), [Roma](#)^c; [Meeus](#), [Mira](#)^{a,d,e}; [Coppieters](#), [Iris](#)^{a,f,d}; [Caeyenberghs](#), [Karen](#)^{a,g}; [Cagnie](#), [Barbara](#)^a

THE IMPORTANCE OF SOCIAL NETWORKS AMONGST REFUGEES RESETTLED THROUGH THE COMMUNITY SPONSORSHIP SCHEME AND THE VULNERABLE PERSONS RESETTLEMENT SCHEME

Critical transmission sectors in embodied atmospheric mercury emission network in China

[Kehan He](#)¹ | [Zhifu Mi](#)¹ | [Long Chen](#)² | [D'Maris Coffman](#)¹ | [Sai Liang](#)³



Fast unfolding of communities in large networks

VD Blondel, JL Guillaume, R Lambiotte... - Journal of statistical ..., 2008 - iopscience.iop.org

We propose a simple method to extract the community structure of large networks. It is a heuristic method that is based on modularity optimization. It is shown to outperform other methods.

☆ Save ⌂ Cite Cited by 25562 Related articles All 35 versions

OPEN Inferring strategies from observations in long iterated Prisoner's dilemma experiments

Eladio Montero-Porras^{1,5}, Jelena Grujic^{1,5}, Elias Fernández Domingos^{1,2} & Tom Lenaerts^{1,2,3,4}

Cell Reports

Article

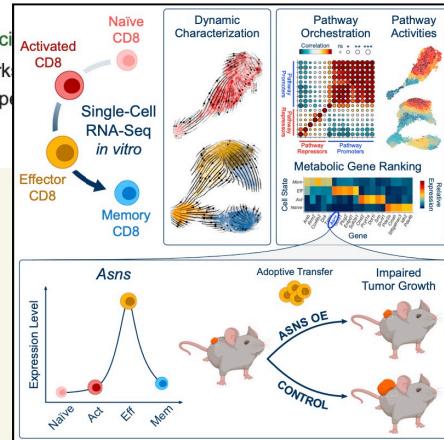
CD8⁺ T cell metabolic rewiring defined by scRNA-seq identifies a critical role of ASNS expression dynamics in T cell differentiation

Untangling the
STRUCTURE and DYNAMICS
of ecological networks

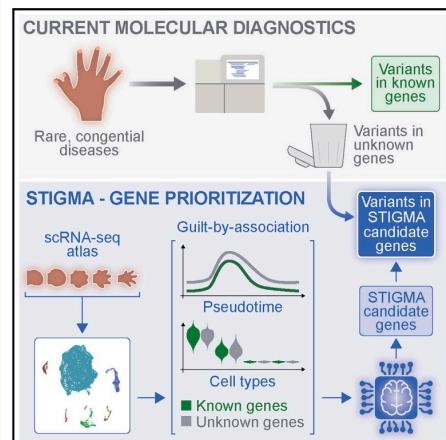
Bernat Bramon Mora

June 1, 2019

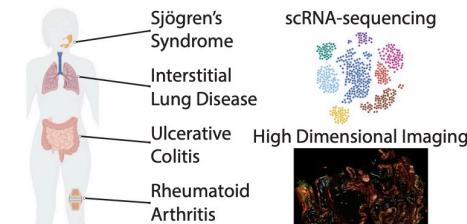
Graphical abstract



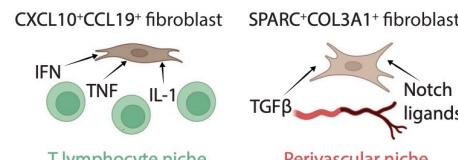
Graphical abstract



Profiling fibroblasts in inflammation disease



Niche signals drive inflammatory phenotypes



Modularity Louvain alg.

- set $A = \{\{v_1\}, \{v_2\}, \dots, \{v_n\}\}$

- (I) • pick a unif. random labelling of vertices 1, ..., n

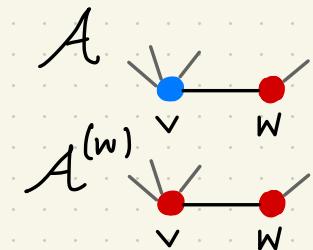
- for $v \in 1, \dots, n$

- for each w nbr of v

- construct $A^{(w)}$ re-colour v with colour of w.

- if $q_{A^{(w)}} > q_A$ $A \rightarrow A^{(w)}$

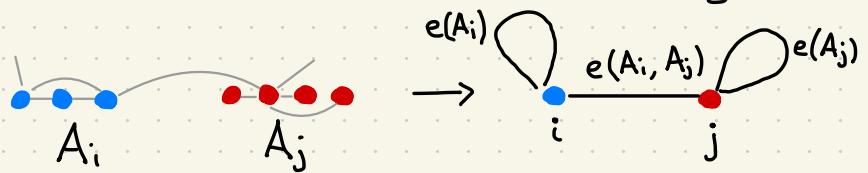
- if no change \rightarrow (II)



- (II) • construct G' by shrinking each colour class to a single vertex

- if $G' = G$ output A

else \rightarrow (I)



Modularity Louvain alg.

- set $A = \{\{v_1\}, \{v_2\}, \dots, \{v_n\}\}$

- ① • pick a unif. random labelling of vertices 1, ..., n

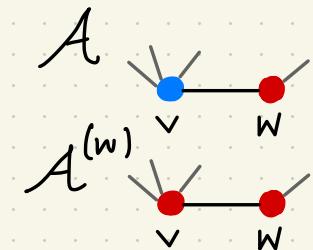
- for $v \in 1, \dots, n$

- for each w nbr of v

- construct $A^{(w)}$ re-colour v with colour of w.

- if $q_{A^{(w)}} > q_A$ $A \rightarrow A^{(w)}$

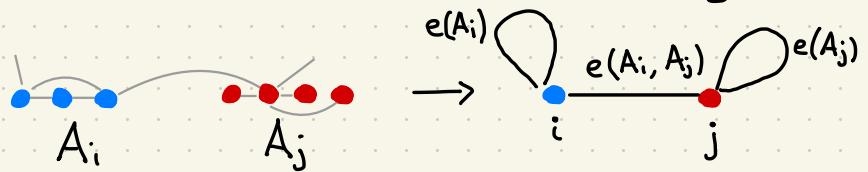
- if no change \rightarrow ②



- ② • construct G' by shrinking each colour class to a single vertex

- if $G' = G$ output A

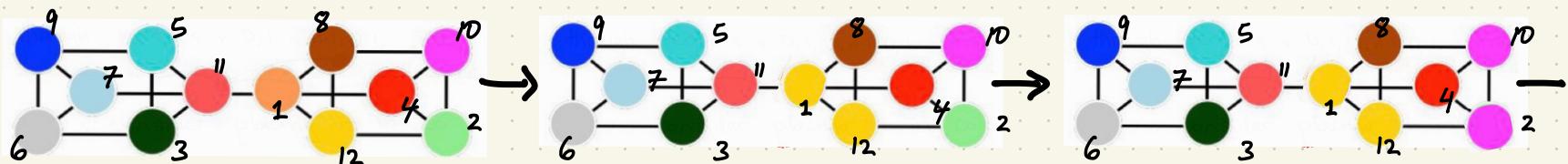
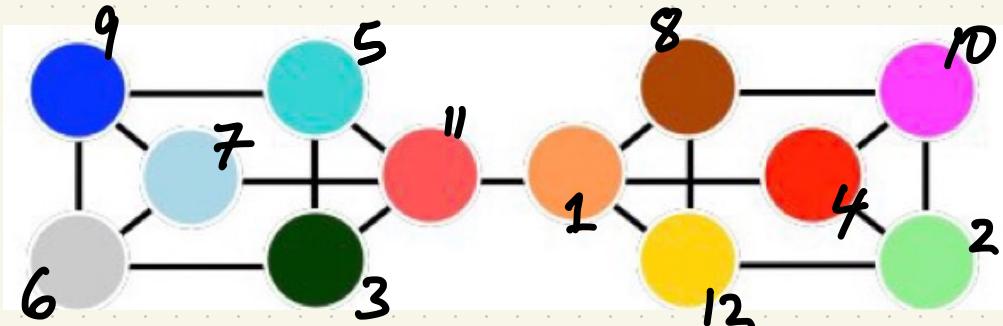
else \rightarrow ①



weighted graph, with loops.

Modularity Louvain alg.

- pick a unif. random labelling of vertices $1, \dots, n$

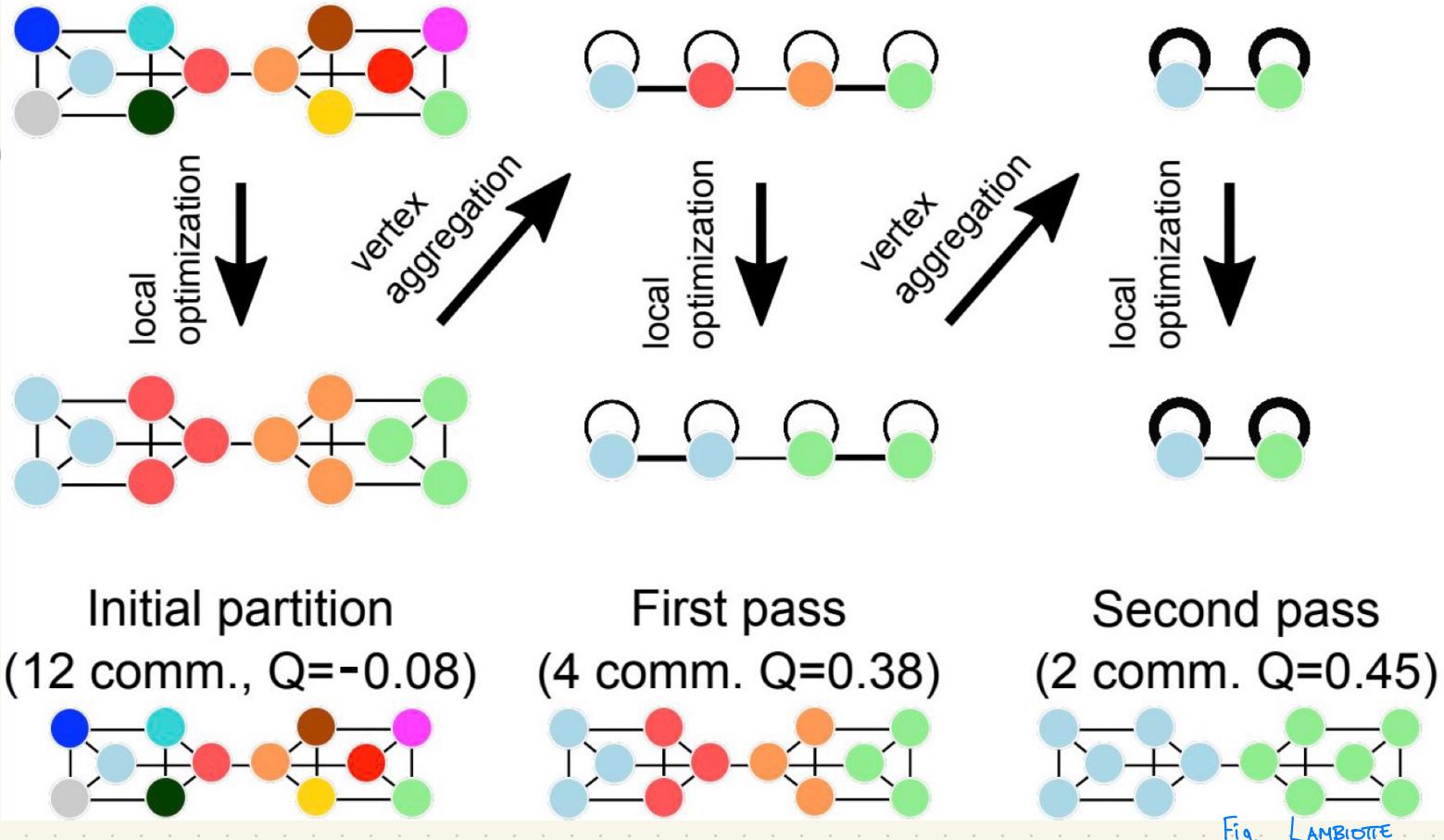


etc; until no local move increases modularity ...

Fig. LAMBIOTTE

Modularity

Louvain alg.



Modularity Louvain alg.

- set $A = \{\{v_1\}, \{v_2\}, \dots, \{v_n\}\}$

- (I) • pick a unif. random labelling of vertices 1, ..., n

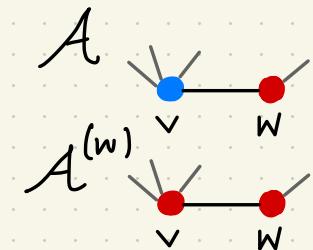
- for $v \in 1, \dots, n$

- for each w nbr of v

- construct $A^{(w)}$ re-colour v with colour of w.

- if $q_{A^{(w)}} > q_A$ $A \rightarrow A^{(w)}$

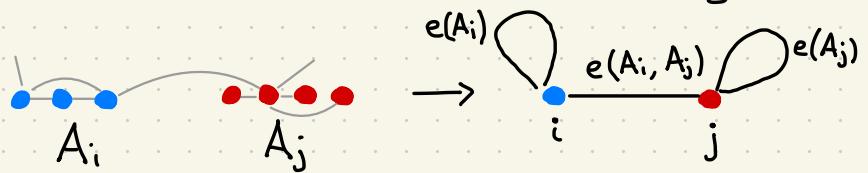
- if no change \rightarrow (II)



- (II) • construct G' by shrinking each colour class to a single vertex

- if $G' = G$ output A

else \rightarrow (I)



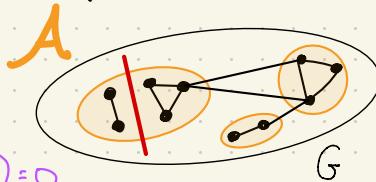
weighted graph, with loops.

Modularity Properties

$A \in \text{OPT}(G)$ i.e. $q_A(G) = q^*(G)$

$\Rightarrow \forall A \in A \quad G[A]$ conn. (+ isolated vert)

\Rightarrow pendant vertex
in same part $\Rightarrow q^*(\text{star}) = 0$



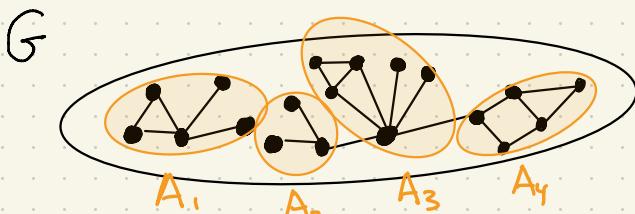
Modularity value:

Robust to small perturbations in edge set

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

$$\forall \lambda: |q_\lambda(G) - q_\lambda(G \setminus E)| < \frac{2|E|}{e(G)}$$

Modularity 'meas of how well a graph can be clustered'

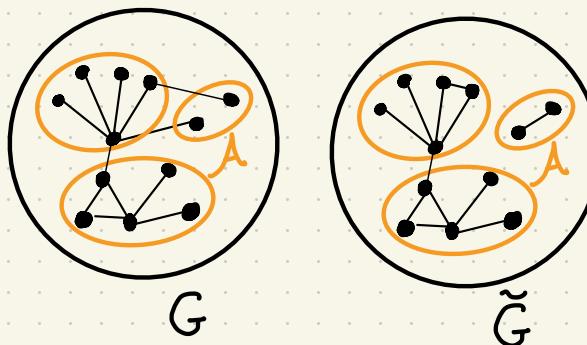
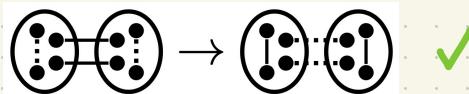


$$A = \{A_1, \dots, A_4\}$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbb{1}_{[u \sim v]} - \frac{d_u d_v}{2m}$$

"edge contrib." "degree tax"

Fix A , which $G \rightarrow \tilde{G}$ ensures $q_A(\tilde{G}) > q_A(G)$?



graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

$$q^*(G) = \max_A q_A(G)$$

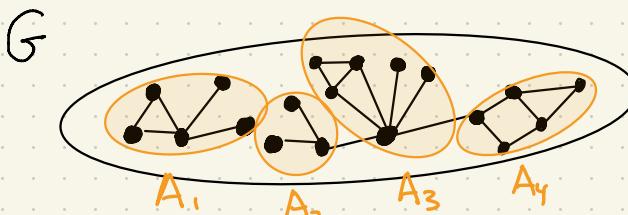
"high vals taken to indicate
more community structure"

d_u = # edges incident to u

$e(A)$ = # edges in set A

$$\text{vol}(A) = \sum_{u \in A} d_u$$

Modularity 'meas of how well a graph can be clustered'



$$A = \{A_1, \dots, A_4\}$$

graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G

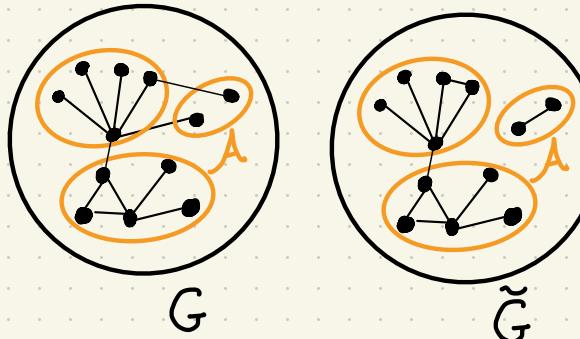
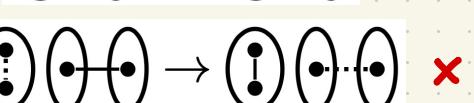
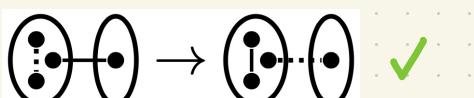
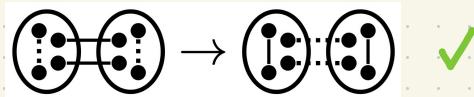
$$q^*(G) = \max_A q_A(G)$$

"high vals taken to indicate more community structure"

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2} = \frac{1}{2m} \sum_{A \in A} \sum_{u, v \in A} \mathbf{1}_{[u \sim v]} - \frac{d_u d_v}{2m}$$

"edge contrib." "degree tax"

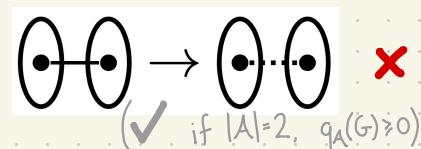
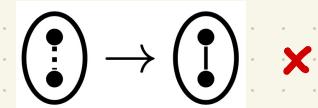
Fix A , which $G \rightarrow \tilde{G}$ ensures $q_A(\tilde{G}) > q_A(G)$?



d_u = # edges incident to u

$e(A)$ = # edges in set A

$$\text{vol}(A) = \sum_{u \in A} d_u$$



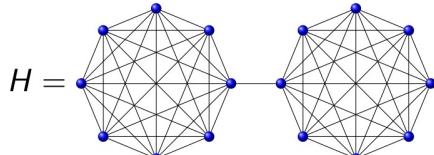
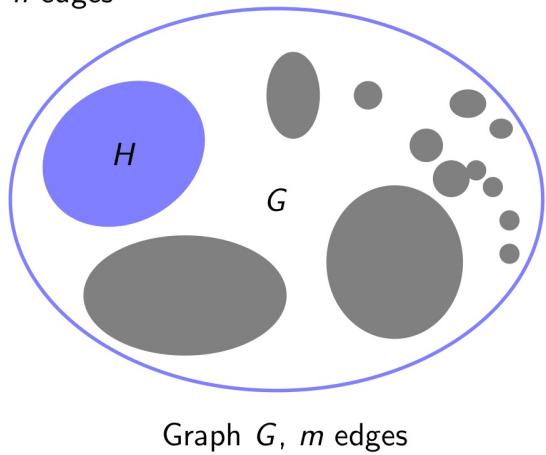
(✓ if $|A|=2$, $q_A(G) \geq 0$)

Modularity Properties

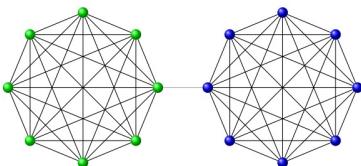
- Resolution limit OPT partition fragile

FORTUNATO AND BARTHÉLEMY 08

Subgraph H
 h edges



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



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

- Modularity value: Robust to small perturbations in edge set

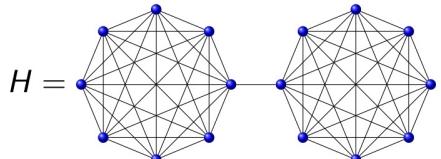
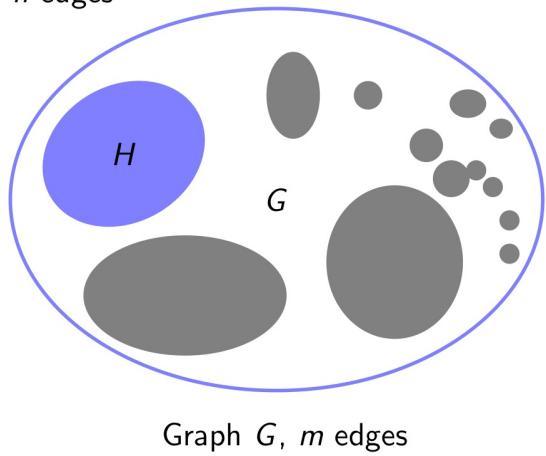
$$\left| q^*(G) - q^*(G \setminus E) \right| < \frac{2|E|}{e(G)}$$

Modularity Properties

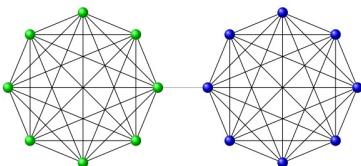
- Resolution limit OPT partition fragile

FORTUNATO AND BARTHÉLEMY 08

Subgraph H
 h edges



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

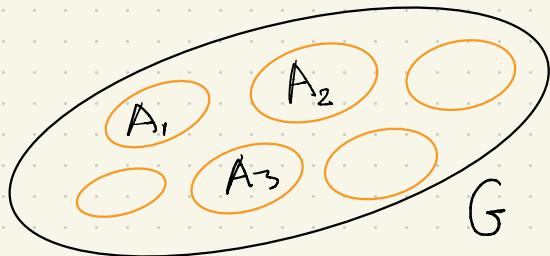


If $h > \sqrt{2m}$, e.g. $m = 1624$.

- Modularity value: Robust to small perturbations in edge set

$$\forall \Delta: |q_L(G) - q_L(G \setminus E)|, |q^*(G) - q^*(G \setminus E)| < \frac{2|E|}{e(G)}$$

Modularity 'meas. of how well a graph can be clustered'



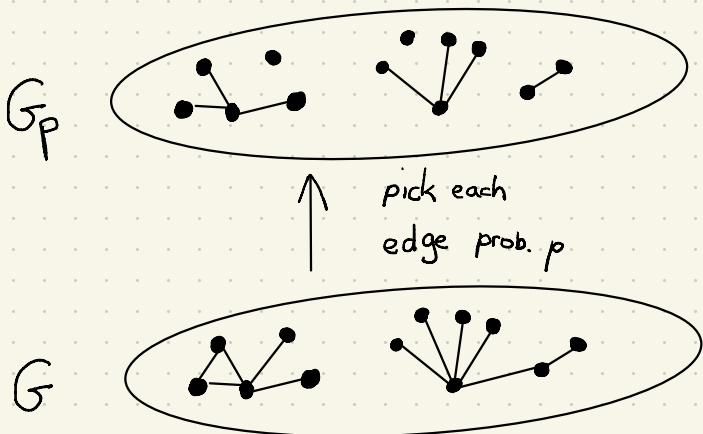
graph G , m edges. $A = \{A_1, \dots, A_k\}$ vertex partition

score of partition A , $q_A(G) =$

modularity of G $q^*(G) = \max_A q_A(G)$

"high vals taken to indicate
more community structure"

Sampling



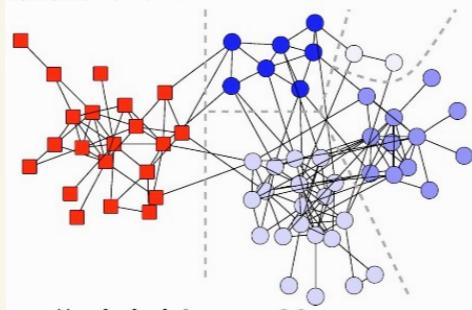
$$G_p = (V, E_p)$$

E_p each edge
kept indep. prob. p

$$G = (V, E) \text{ fixed graph}$$

Simulations

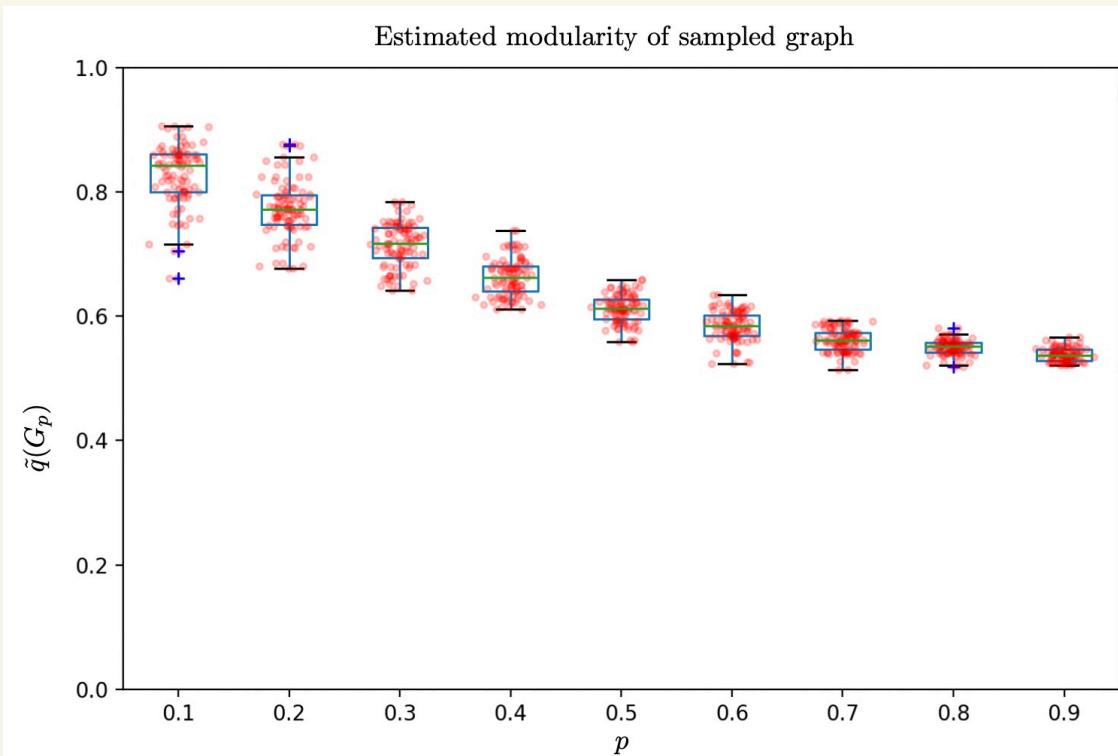
Dolphin Network [Lusseau]



$$|V| = 62 \quad |E| = 152$$

$$q^*(G) = 0.529\dots \text{ (3 dec places)}$$

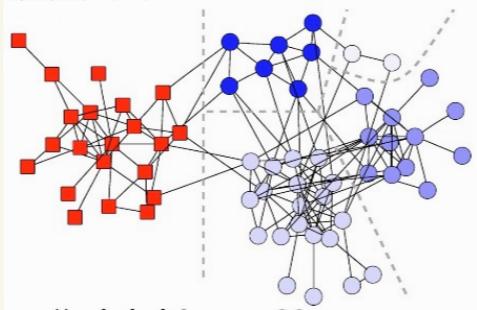
[BRANDES + '08]



To estimate modularity take max of 200 runs of Lovain and Leiden algs.

Simulations

Dolphin Network [Lusseau]

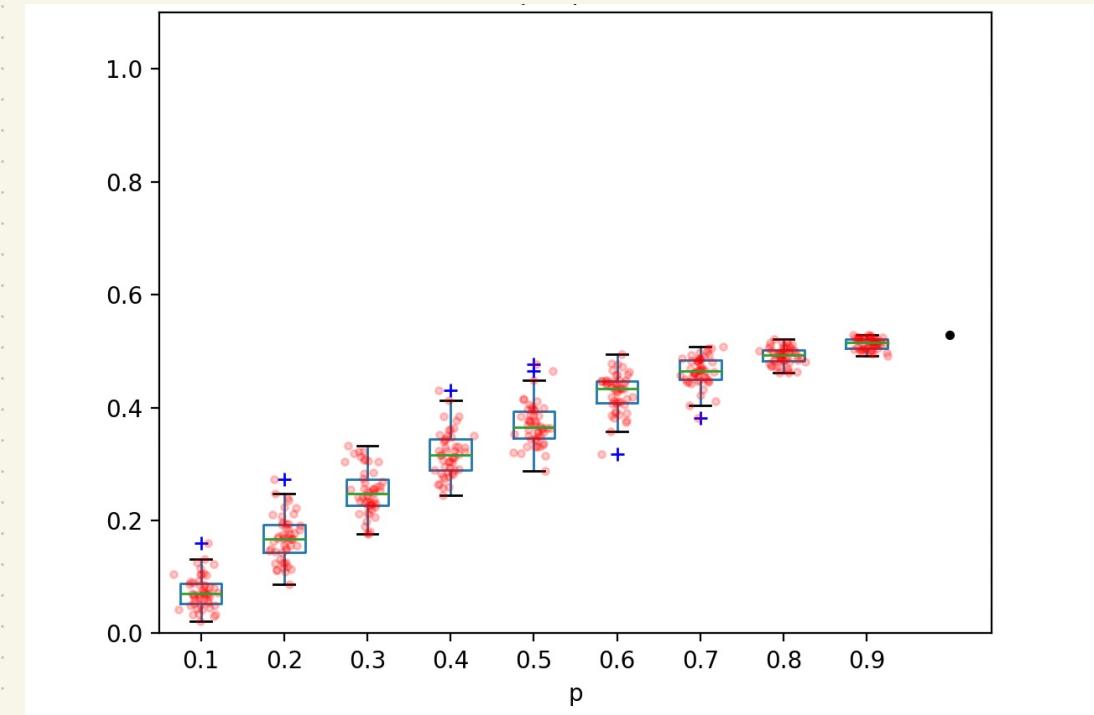


$$|V| = 62 \quad |E| = 152$$

$$q^*(G) = 0.529\dots \text{ (3 dec places)}$$

[BRANDES + '08]

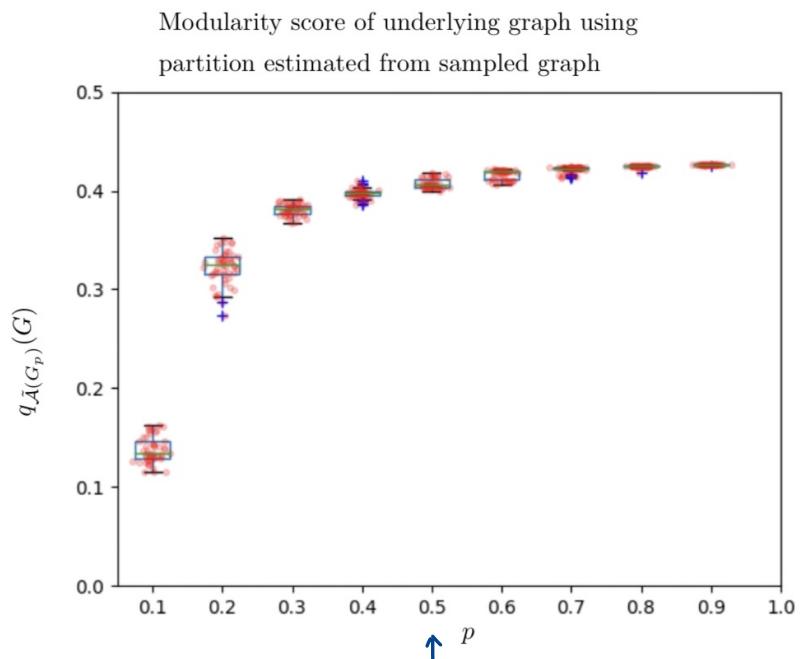
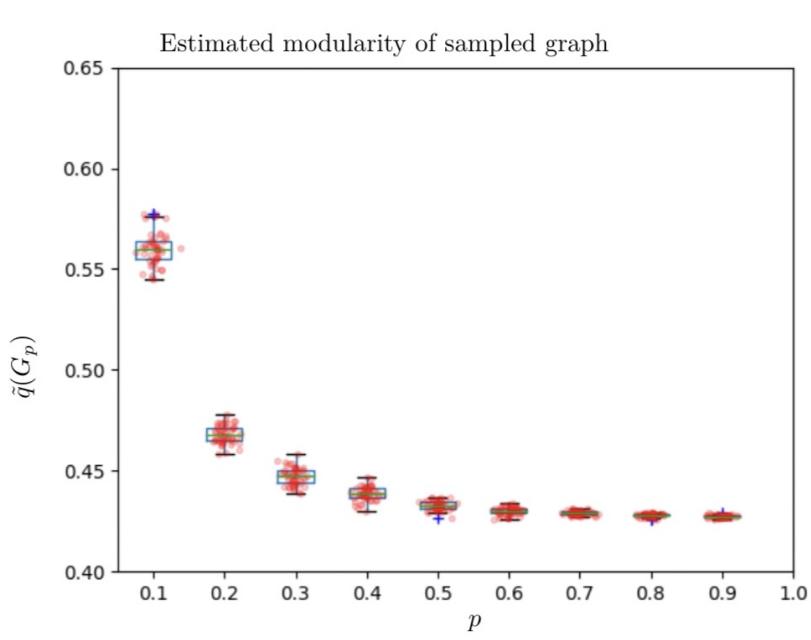
Modularity score of underlying dolphin network using partition $\sim \text{OPT}$ of sampled network.



To estimate OPT of sampled G_p take max of 200 runs of Lovain and Leiden algs.
↑
on G_p

US Political Blogs [ADAMIC, GLANCE '2005]

Graph $V \sim 1500$ $E \sim 16000$



seeing half the edges
≈ " all " "]
'how good
the partition is'

II: Fundamental limits of learning

- when can we detect / recover planted communities ?
- when can we do this fast ?

Planted Community $G \sim G(n, p, q^*, K)$, $i \in K$ w prob $\frac{K}{n}$

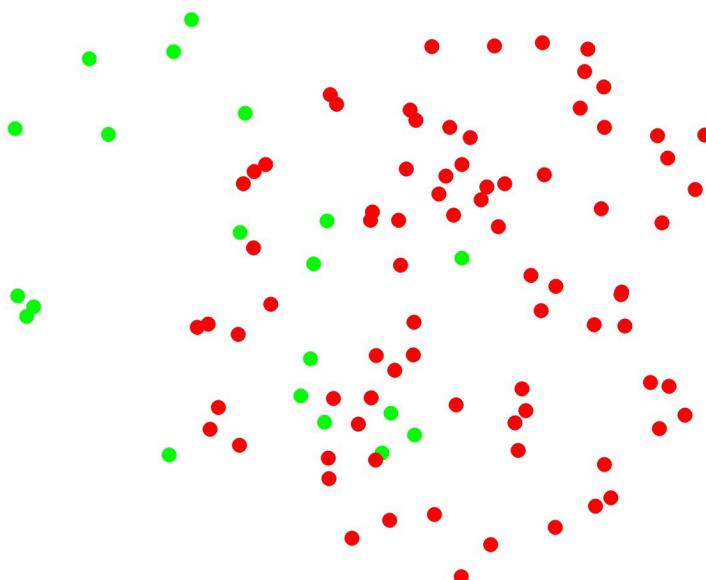
$\overset{\text{signal}}{\downarrow}$ $\overset{\text{noise}}{\downarrow}$

$p > q$

$$A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{ow} \end{cases}$$

n points

- $\sim K$ 'community' nodes
- $\sim n-K$ 'non-community' "



Planted Community $G \sim G(n, p, q^*, K)$, $i \in K$ w prob $\frac{K}{n}$

$$A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{otherwise} \end{cases}$$

$$p > q$$

n points

- K 'community' nodes
- $n-K$ 'non-community' "
- with prob. P

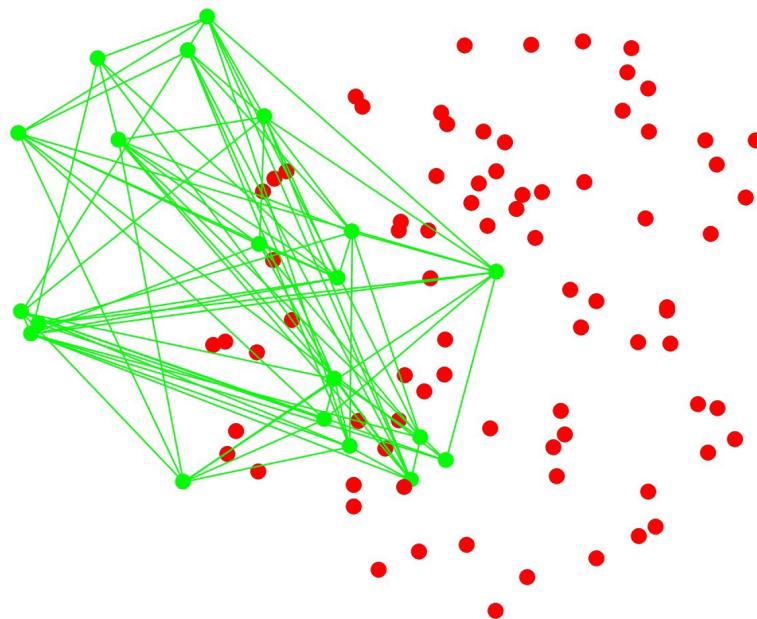


Fig: Jianming Xu, Duke

Planted Community

$$G \sim G(n, p, q', k), \quad i \in K \text{ w.prob } \frac{k}{n} \quad A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{otherwise} \end{cases}$$

↑
 $p > q'$

n points

- K 'community' nodes
- n-k 'non-community' "

- —● with prob. P
- ● " " q
- ● " " q

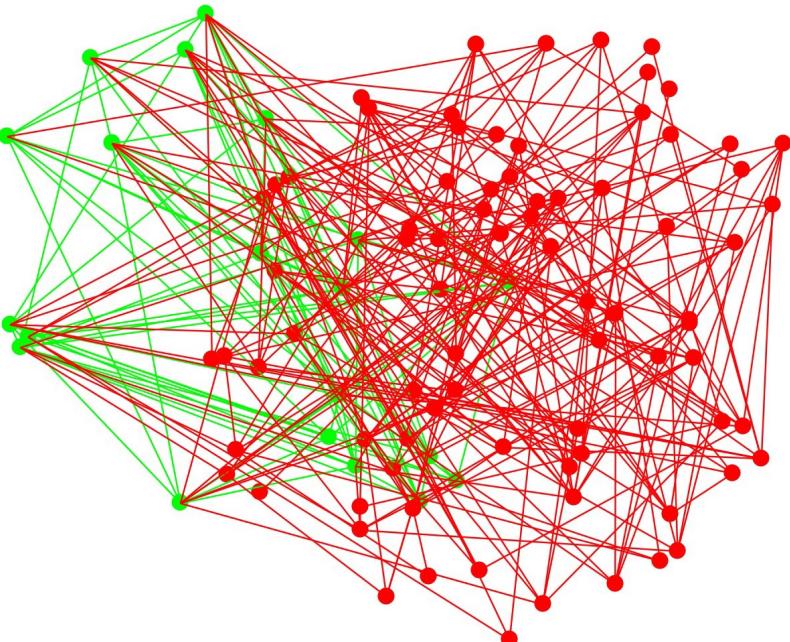


Fig: Jianming Xu, Duke

Planted Community $G \sim G(n, p, q^*, K)$, $i \in K$ w.prob $\frac{k}{n}$ $A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{otherwise} \end{cases}$

Process

n points

- K 'community' nodes
- $n-K$ 'non-community' "

with prob. p

" " q

" " q

Output

unlabelled graph

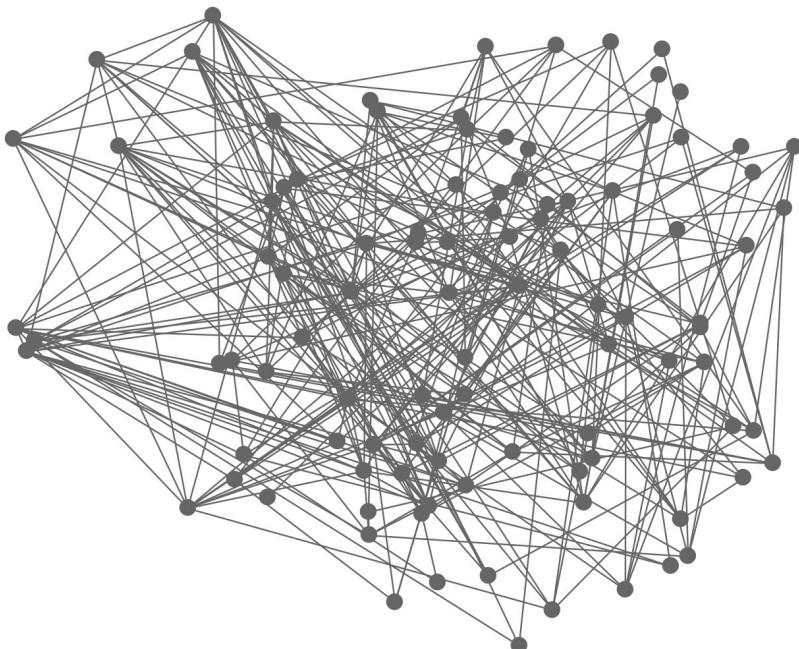


Fig: Jianming Xu, Duke

Planted Community

$$G \sim G(n, p, q^*, K), \quad i \in K \text{ w prob } \frac{k}{n}$$

↑
 $p > q$

$$A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{ow} \end{cases}$$

Process

n points

- K 'community' nodes
- n-k 'non-community' "

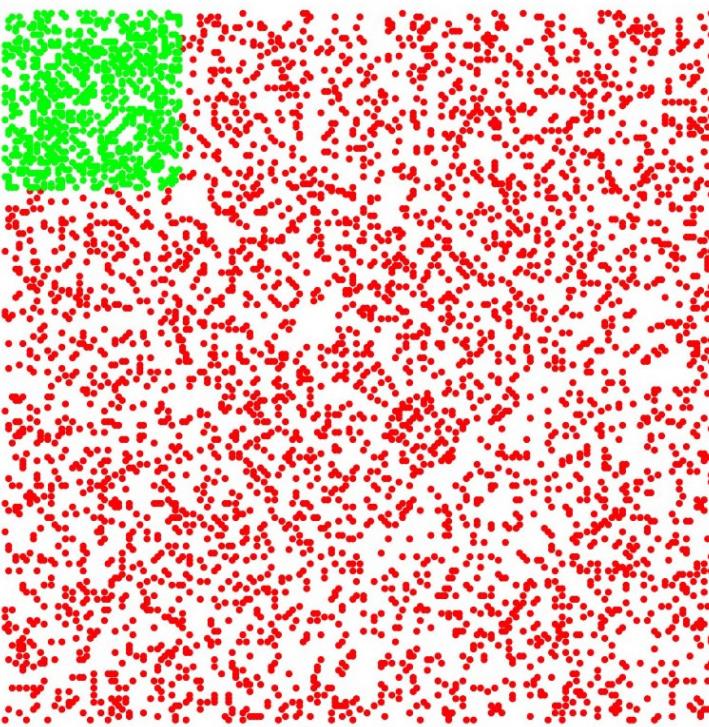
with prob. P

" " q

" " q

Output

unlabelled graph



$n=200$

$K=50$

$p=0.3$

$q=0.1$

Fig: Jianming Xu, Duke

Planted Community

$$G \sim G(n, p, q^*, K), \quad i \in K \text{ w prob } \frac{k}{n} \quad A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{ow} \end{cases}$$

↑
signal noise
 $p > q$

Process

n points

- K 'community' nodes
- $n-K$ 'non-community' "

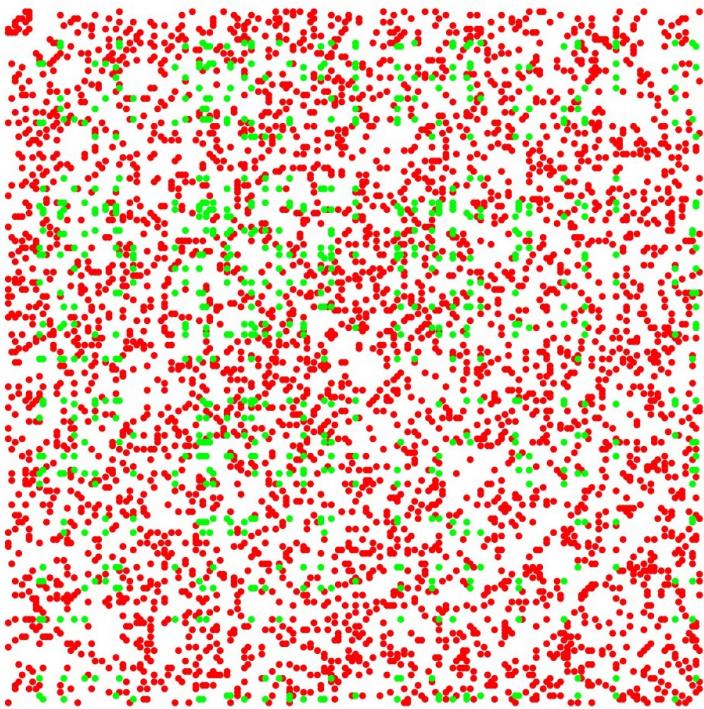
with prob. P

" " q

" " q

Output

unlabelled graph



$n=200$

$k=50$

$p=0.3$

$q=0.1$

Fig: Jianming Xu, Duke

Planted Community

$$G \sim G(n, p, q^*, K), \quad i \in K \text{ w.prob } \frac{k}{n} \quad A_{ij} = \begin{cases} Be(p) & i, j \in K \\ Be(q) & \text{otherwise} \end{cases}$$

↑
 $p > q^*$

Process

n points

- K 'community' nodes
- $n-K$ 'non-community' "

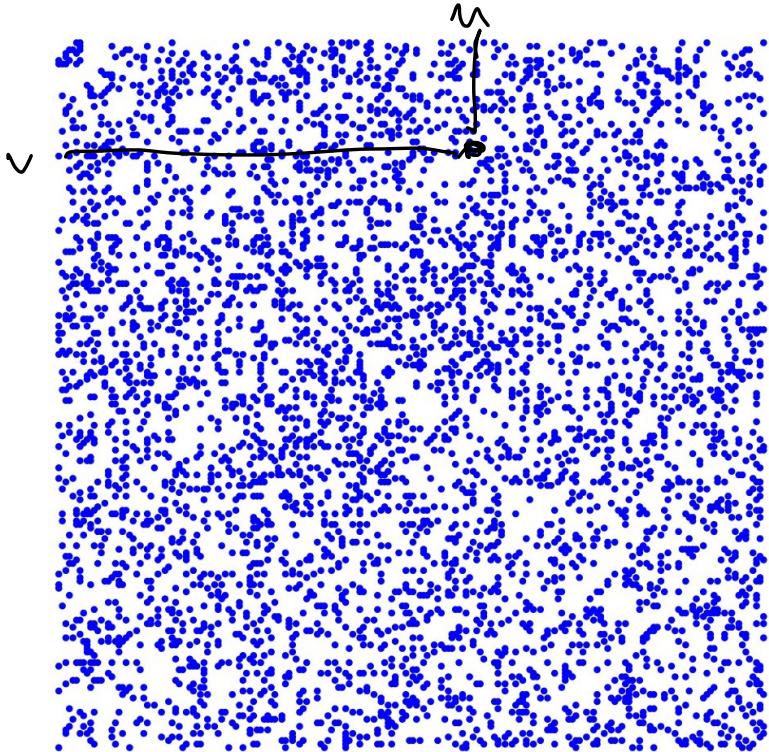
 with prob. P

 " " q^*

 " " " q

Output

unlabelled graph



$n=200$

$K=50$

$p=0.3$

$q=0.1$

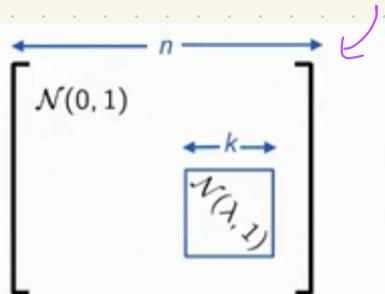
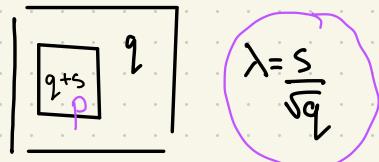
Fig: Jianming Xu, Duke

Planted Dense subgraph / Submatrix

Hypothesis Testing : asymptotics

$$H_0: G(n, q)$$

$$H_1: G(n, k, s, q)$$



$$H_0: \text{i.i.d. } N(0, 1)$$

$$H_1: \text{submatrix of } N(\lambda, 1)$$

$$\bullet \quad n \rightarrow \infty$$

- for what - $k(n)$ **size** of planted structure
- $\lambda(n)$ **strength** of signal
- can we find test ϕ

$$P_0(\underline{\phi(Y) = 1}) + P_1(\underline{\phi(Y) = 0}) \rightarrow 0$$

- when is there a **fast test** ?

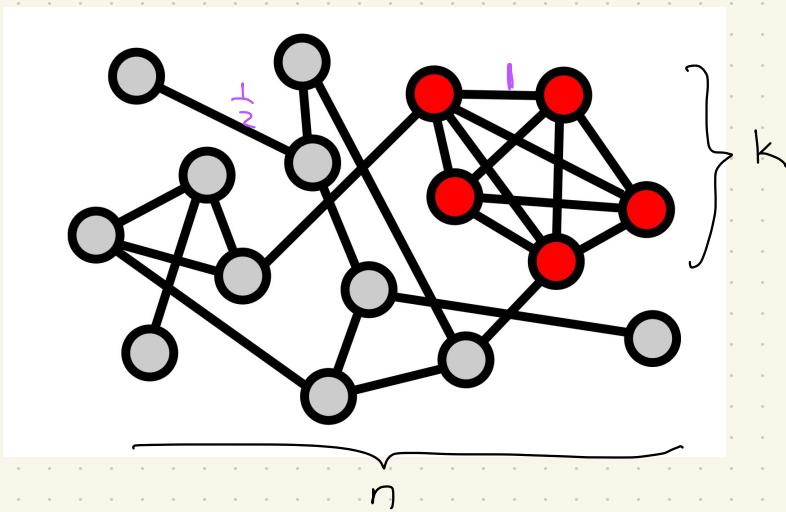
Planted Clique $G \sim G(n, \frac{1}{2}, k)$, $v \in K$ with prob. $\frac{k}{n}$

$$A_{uv} = \begin{cases} 1 & u, v \in K \\ \text{Be}(\frac{1}{2}) & \text{ow} \end{cases}$$

Two parameters

- size of planted structure
- size of entire network

Q: When can we find planted clique?



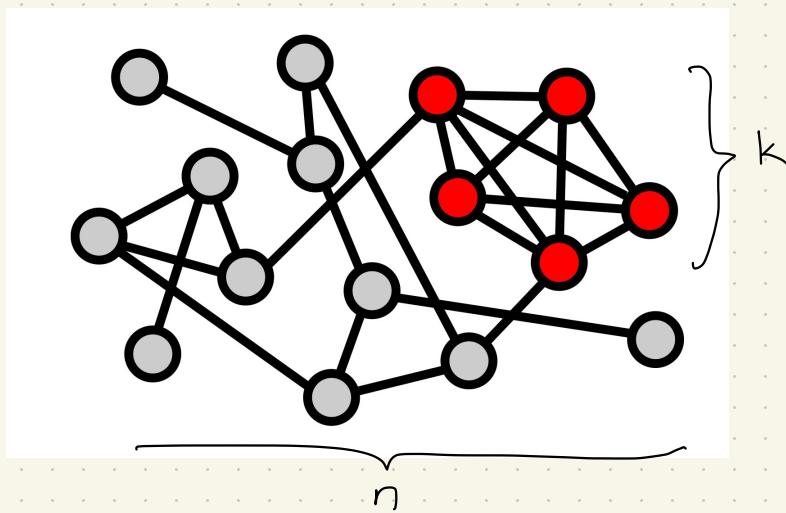
Planted Clique $G \sim G(n, \frac{1}{2}, k)$, $v \in K$
 with prob. $\frac{k}{n}$

$$A_{uv} = \begin{cases} 1 & u, v \in K \\ Be\left(\frac{1}{2}\right) & \text{ow} \end{cases}$$

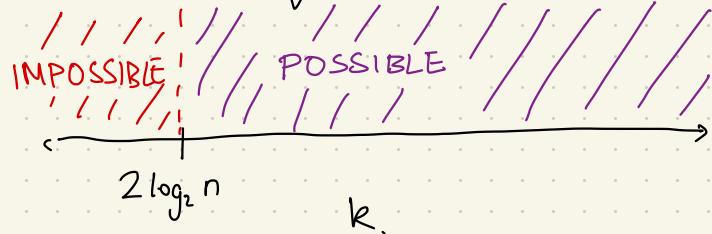
~~IMPOSSIBLE~~

$2 \log_2 n$ k

$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one
 in amongst 'background' one.



Planted Clique $G \sim G(n, \frac{1}{2}, k)$, $\forall v \in K$ with prob. $\frac{k}{n}$

$$A_{uv} = \begin{cases} 1 & u, v \in K \\ Be\left(\frac{1}{2}\right) & \text{otherwise} \end{cases}$$


$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

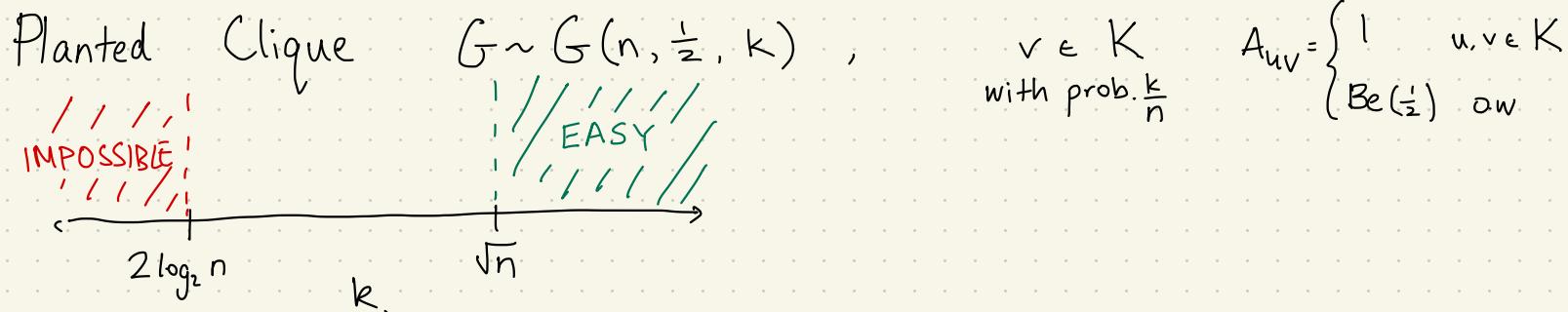
Methods to find clique

① BRUTE-FORCE

Search all k -vertex subsets

\hat{K} first clique found.

if $K \geq (2+\varepsilon) \log n$ $P(\hat{K} = K) \rightarrow 1$



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

Methods to find clique

① DEGREE TEST

\hat{K} = set of K vertices of highest degree

$$K \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

Planted Clique $G \sim G(n, \frac{1}{2}, k)$, $v \in K$ with prob. $\frac{k}{n}$ $A_{uv} = \begin{cases} 1 & u, v \in K \\ \text{Be}(\frac{1}{2}) & \text{ow} \end{cases}$

IMPOSSIBLE!

$2 \log_2 n$ k .

$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one in amongst 'background' one.

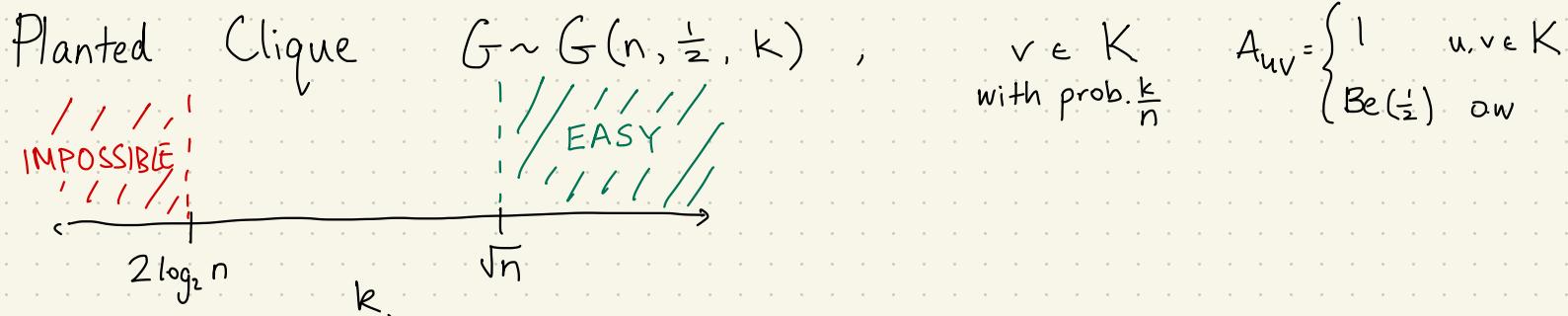
Methods to find clique

① DEGREE TEST

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n \log n} \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

$$K \in \binom{[n]}{k}$$



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

Methods to find clique

① DEGREE TEST

\hat{K} = set of K vertices of highest degree

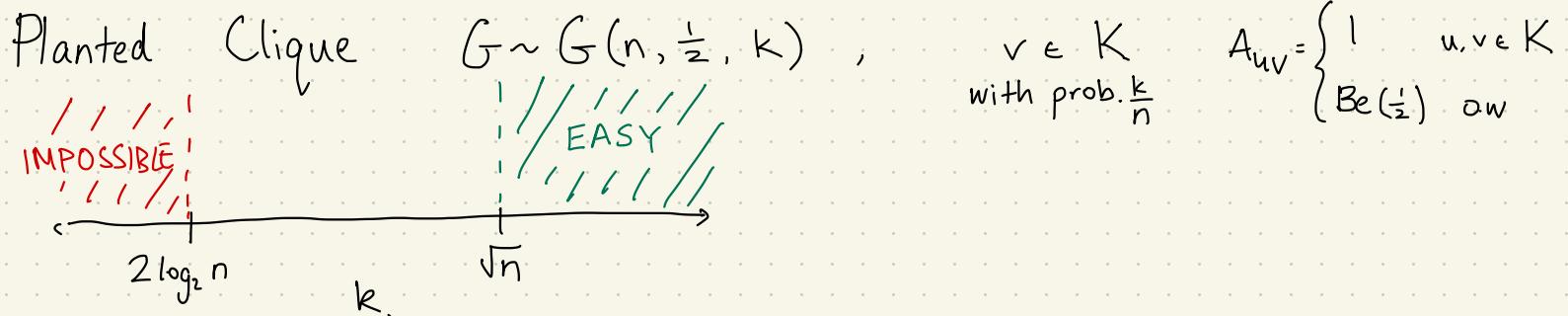
$$K \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

② SPECTRAL METHOD

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of K largest $|u_i|$



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

Methods to find clique

① DEGREE TEST

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

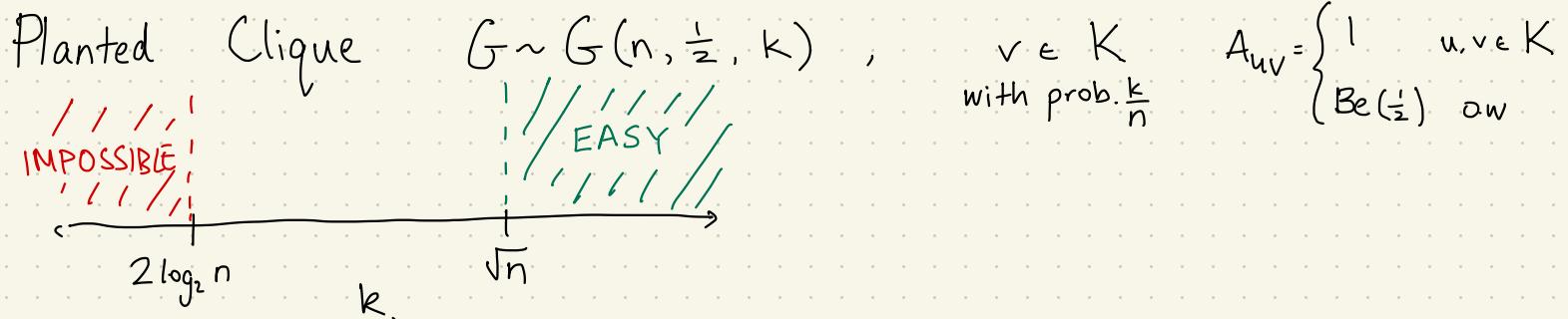
② SPECTRAL METHOD

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of k largest $|u_i|$

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1$$



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

Methods to find clique

① **DEGREE TEST**

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

② **SPECTRAL METHOD**

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

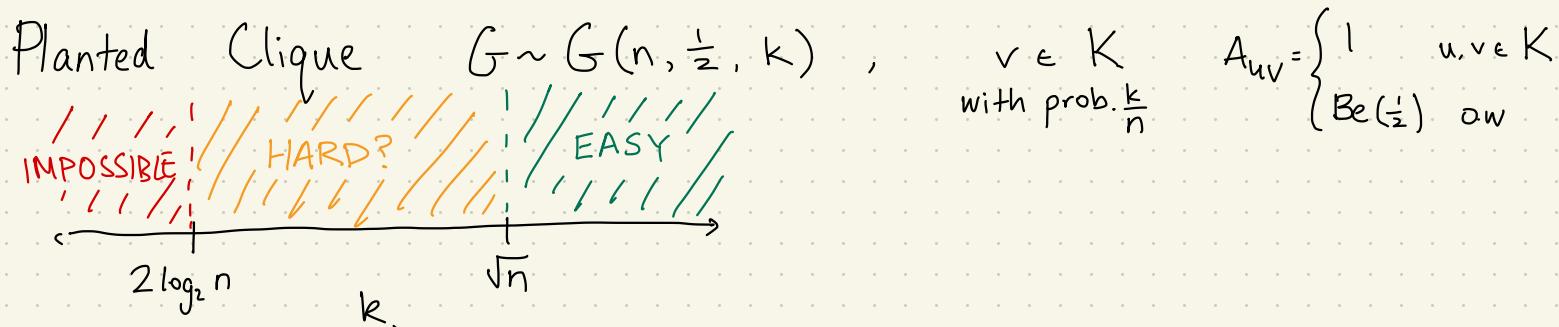
(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of k largest $|u_i|$

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1$$

③ **SDP METHOD**

Yes. If $k = \mathcal{O}(\sqrt{n})$.



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$. \Rightarrow can't find 'planted' one amongst 'background' one.

Methods to find clique

① DEGREE TEST

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

② SPECTRAL METHOD

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

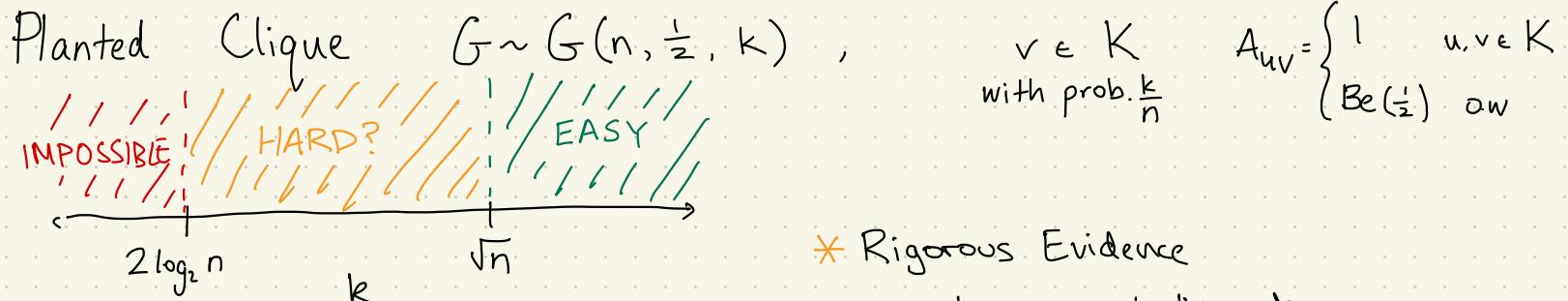
(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of k largest $|u_i|$

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1$$

③ SDP METHOD

Yes. If $k = \mathcal{O}(\sqrt{n})$.



$G' \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$.

Methods to find clique

① DEGREE TEST

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

② SPECTRAL METHOD

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of k largest $|u_i|$

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1$$

③ SDP METHOD

Yes. If $k = \mathcal{O}(\sqrt{n})$.

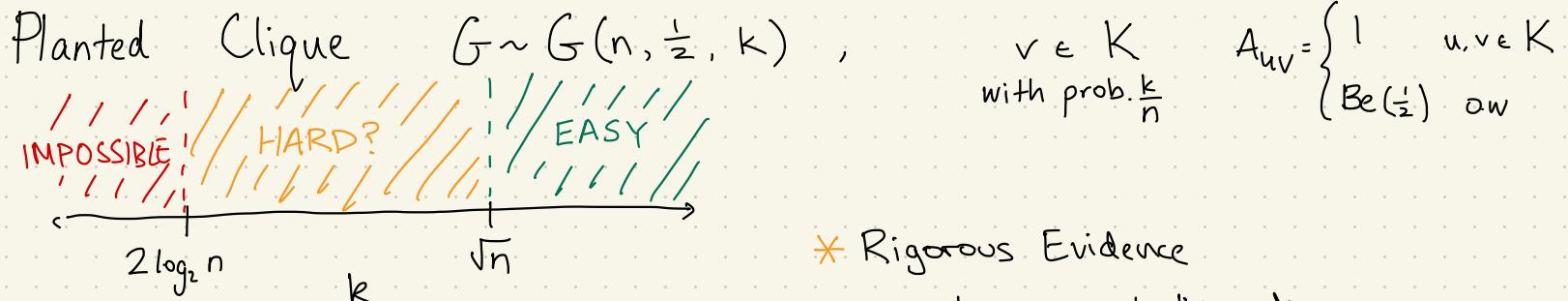
* Rigorous Evidence

suggesting no poly-time alg

- reductions (avg case)
- restricted class of alg

low deg poly

- subgraph tests
- # edges, # A's, ...
- spectral methods.



$G \sim G(n, \frac{1}{2})$: largest clique whp $\sim 2 \log_2 n$.

Methods to find clique

① DEGREE TEST

\hat{K} = set of k vertices of highest degree

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1.$$

② SPECTRAL METHOD

$$W_{ij} = \begin{cases} 2A_{ij} - 1 & i \neq j \\ 0 & \text{o.w} \end{cases}$$

(i) u top eigenvector of W

(ii) (threshold) \hat{K} index vector of k largest $|u_i|$

$$k \geq \sqrt{n} \log n \Rightarrow P(\hat{K} = K) \rightarrow 1$$

③ SDP METHOD

Yes. If $k = \mathcal{O}(\sqrt{n})$.

* Rigorous Evidence

suggesting no poly-time alg

- reductions (avg case)
- restricted class of alg

low deg poly

- subgraph tests
- # edges, # A's, ...
- spectral methods.

PLANTED DENSE SUBGRAPH

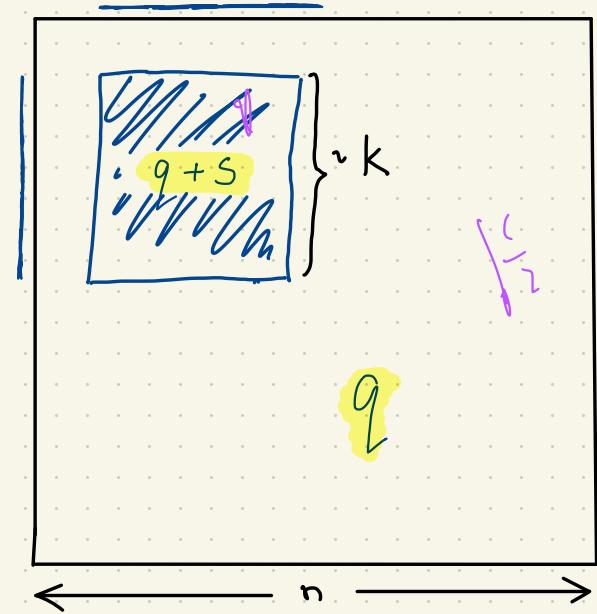
Vertex labels : $\sigma_v = \begin{cases} 1 & \bullet \text{ w. prob } \frac{k}{n} \\ \emptyset & \bullet \text{ w. prob } 1 - \frac{k}{n} \end{cases}$

Observe $Y_{uv} \sim \begin{cases} \text{Ber}(q+s) & \sigma_u = \sigma_v = 1 \\ \text{Ber}(q) & \text{o.w.} \end{cases}$

ALGORITHMIC QNS

- Detection : determine if whp sample from planted model or not
- Recovery : given sample from planted model find community (exactly? weakly corr?)
- "Counting" ... ?

$$G(n, k, q, s)$$



CONTEXT

PLANTED DENSE SUBGRAPH

Detection

$$k = \Theta(n^\beta)$$

EASY

β

1

$\frac{1}{2}$

HARD

IMPOSSIBLE

bigger
planted
structure

decreasing signal

EASY

Recovery

$$k = \Theta(n^\beta)$$

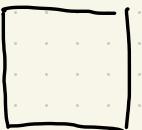
1

$\frac{1}{2}$

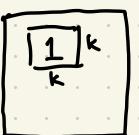
HARD

IMPOSSIBLE

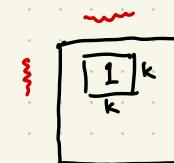
H_0



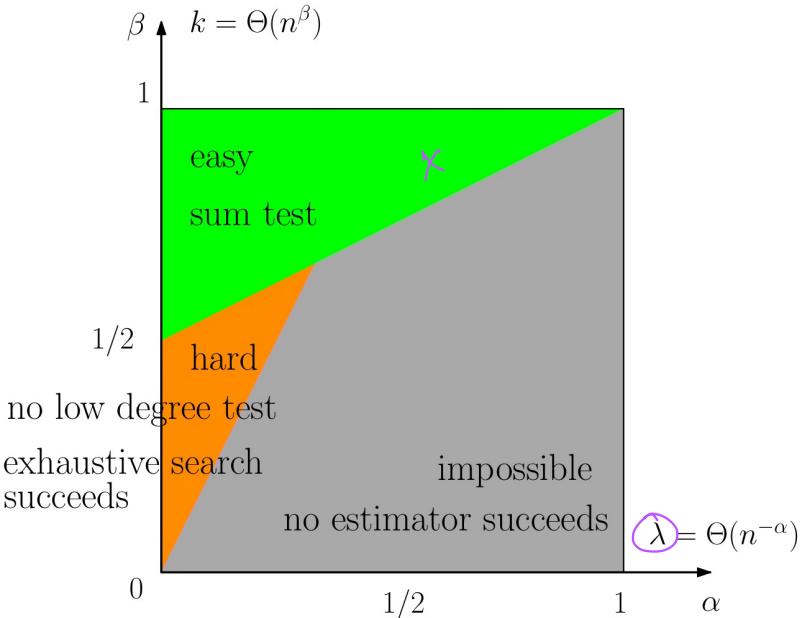
vs. H_1



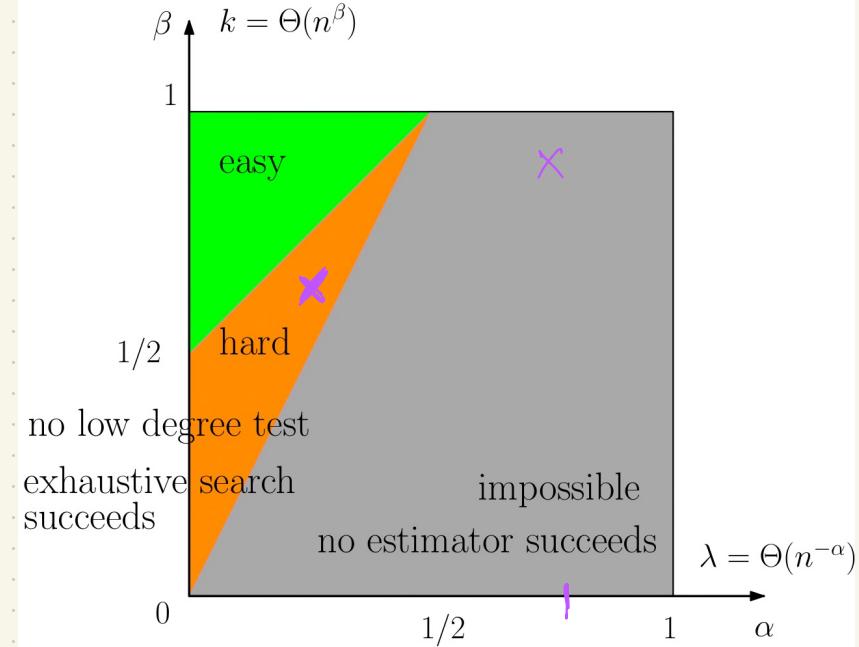
recover



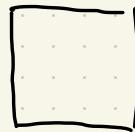
Detection



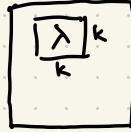
Recovery



H_0



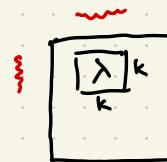
vs. H_1



$$\mathcal{N}(\lambda, 1)$$

$$\lambda = \frac{s}{\sqrt{q}}$$

recover

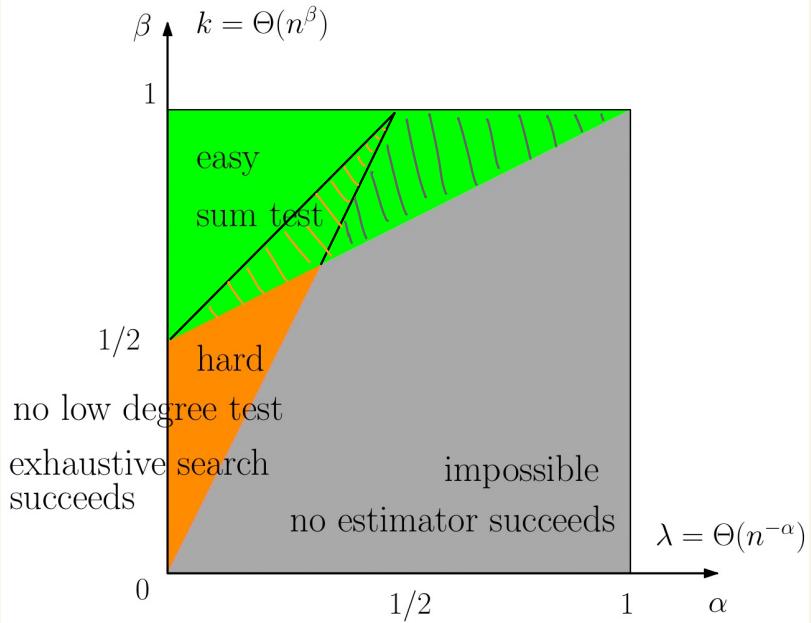


REFS: MANY AUTHORS. BI13, BIS15, MW15, CX16, DM14, CLR17, HWX17, BBH18, GJS19, BMR20, BBP05, BS06, FP07, CDF09, BGN11, SWPN09, KBRS11, BKR⁺11, ACD11, BWZ20, SW22

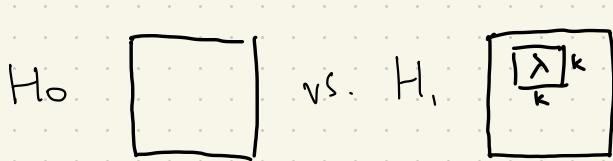
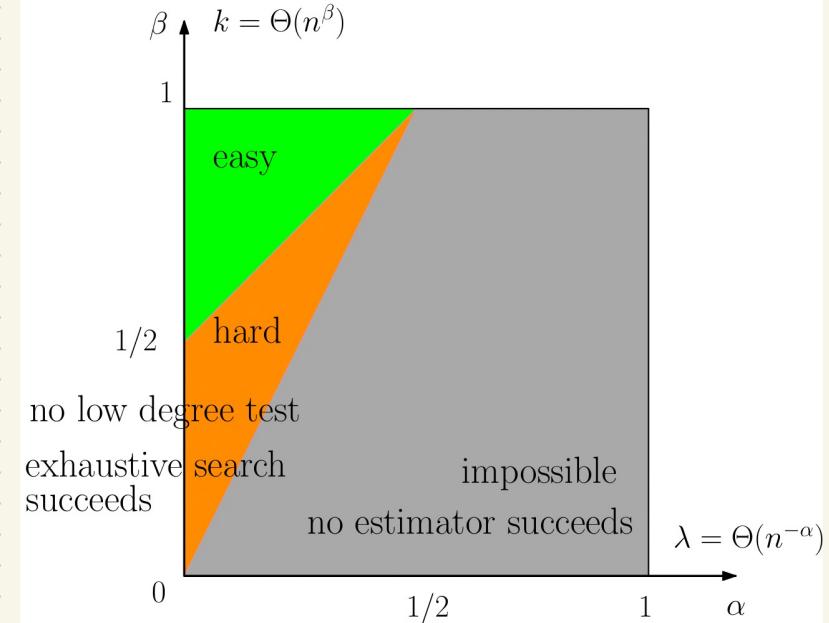
CONTEXT

Detection

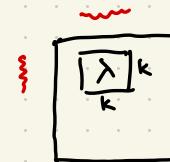
'Easier to detect than recover'.



Recovery

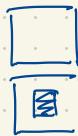


recover



Hypothesis Testing

$$H_0: G \sim P_n = G(n, \frac{1}{2})$$

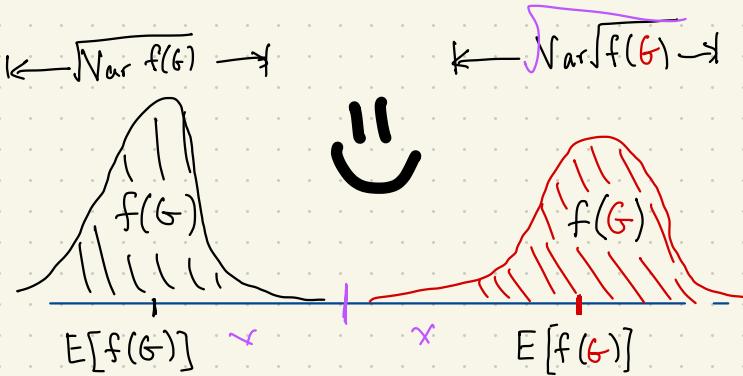


distributions on $\{0,1\}^{(2)}$ or $\mathbb{R}^{(2)}$

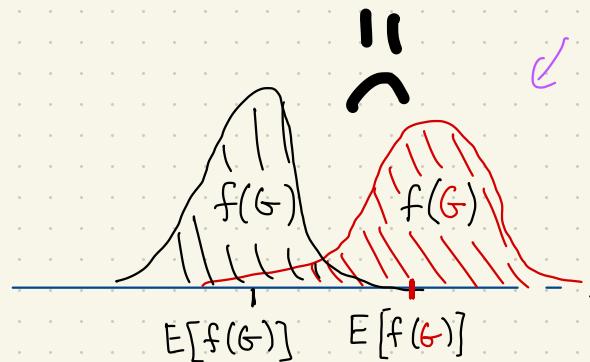
$$H_1: G \sim Q_n = G(n, \frac{1}{2}, k)$$



f detects



f doesn't detect



A 'degree D test' $f_n: \mathbb{R}^{(2)} \rightarrow \mathbb{R}$ $\deg \leq D$.

strongly separates if

$$\mathbb{E}_{P_n}[f] - \mathbb{E}_Q[f] \gg \sqrt{\max\{\text{Var}_Q[f], \text{Var}_P[f]\}}$$

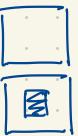
"difference in means" \gg "fluctuations"

Low DEG POLY fail if
no $\Theta(\log(n))$ - degree test.

Hypothesis Testing

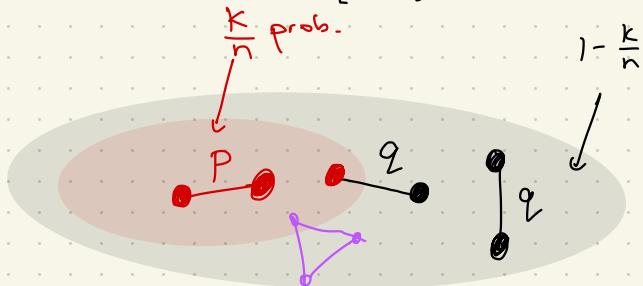
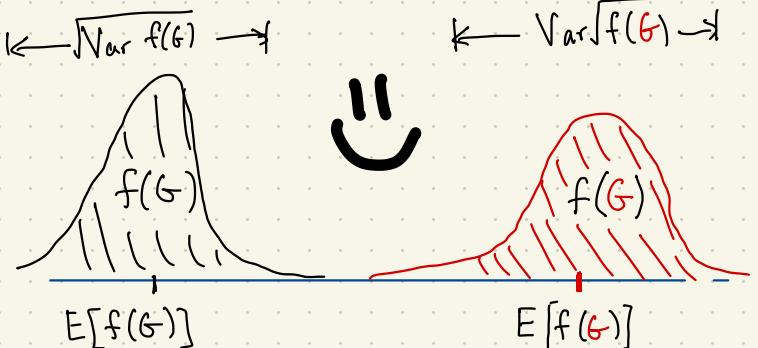
$$H_0: G \sim P_n = G(n, \frac{1}{2})$$

$$H_1: G \sim Q_n = G(n, \frac{1}{2}, k)$$



distributions on
 $\{0,1\}^{\binom{n}{2}}$ or $\mathbb{R}^{\binom{n}{2}}$

f detects



$$E_{P_n}[f] - E_{Q_n}[f] \gg \sqrt{\max\{Var_Q[f], Var_P[f]\}}$$

$$\begin{aligned} \text{Ex: } f &= \sum_{u,v,w} \mathbf{1}_{[uv \in E]} \mathbf{1}_{[uw \in E]} \mathbf{1}_{[vw \in E]} \\ f(G) &= 3! \# \Delta's \end{aligned}$$

calc

$$\begin{aligned} E_{Q_n}[f] &= \sum_{u,v,w} P[uv, uw, vw \in E] q^3 \cdot \left(1 - \frac{k}{n}\right)^3 \\ &= \sum_{u,v,w} P[uv, uw, vw \in E] \begin{bmatrix} u \\ v \\ w \end{bmatrix} \cdot \begin{bmatrix} u \\ v \\ w \end{bmatrix} \\ &\quad + P[uv, uw, vw \in E] \begin{bmatrix} u \\ v \\ w \end{bmatrix} \cdot \begin{bmatrix} u \\ v \\ w \end{bmatrix} \\ &\quad + P[uv, uw, vw \in E] \begin{bmatrix} u \\ v \\ w \end{bmatrix} \cdot \begin{bmatrix} u \\ v \\ w \end{bmatrix} \\ &\quad : \\ &\quad + P[uv, uw, vw \in E] \begin{bmatrix} u \\ v \\ w \end{bmatrix} \cdot \begin{bmatrix} u \\ v \\ w \end{bmatrix} \\ &\quad p^3 \cdot \left(\frac{k}{n}\right)^3 \end{aligned}$$

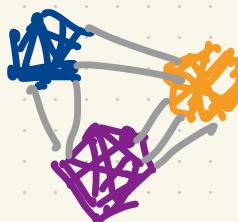
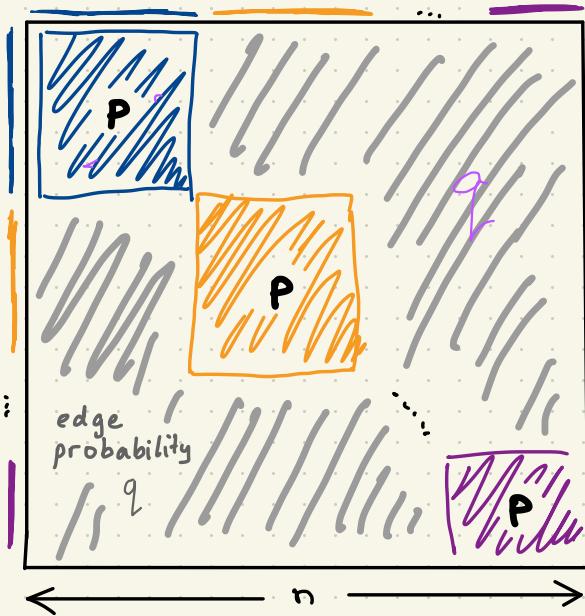
PLANTED PARTITION / STOCHASTIC Block model: M COMMUNITIES

Vertex labels : $\sigma_v = \begin{cases} 1 & \bullet \\ \vdots & \vdots \\ M & \bullet \end{cases}$ w. prob $\frac{1}{M}$

Edges
for $u \sim v$: $A_{uv} \sim \begin{cases} \text{Ber}(p) & \text{if } \sigma_u = \sigma_v \\ \text{Ber}(q) & \text{o.w.} \end{cases}$

ALGORITHMIC QNS

- Detection : determine if whp sample from planted model or not
- Recovery : given sample from planted model find communities (exactly? weakly corr?)
- "Counting" ... ?

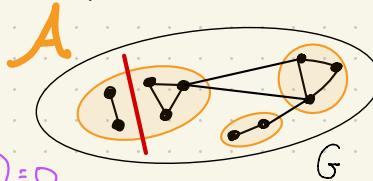


Modularity Properties

$$\lambda \in \text{OPT}(G) \text{ i.e. } q_\lambda(G) = q^*(G)$$

$\Rightarrow \forall A \in \lambda \quad G[A]$ conn. (+ isolated vert)

\Rightarrow pendant vertex in same part $\Rightarrow q^*(\text{star}) = 0$



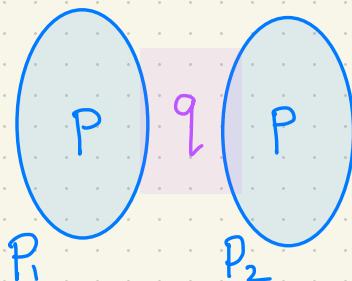
SBM

$$P = \frac{a}{n} \cdot w(n)$$

$$q = \frac{b}{n} \cdot w(n),$$

$$P = \{P_1, P_2\}$$

'planted partition'



Thm [BICKEL HEN] $G \sim G_{n,p,q}$, let $\lambda \in \text{OPT}(G)$

- $w(n) \rightarrow \infty \Rightarrow \text{whp } \lambda \text{ is } o(n) \text{ away } P$
- $\frac{w(n)}{\log n} \rightarrow \infty \Rightarrow " \lambda = P"$

Modularity value:

Robust to small perturbations in edge set

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

$$\forall \lambda: |q_\lambda(G) - q_\lambda(G \setminus E)| < \frac{2|E|}{e(G)}$$

Percolated random graph G_p

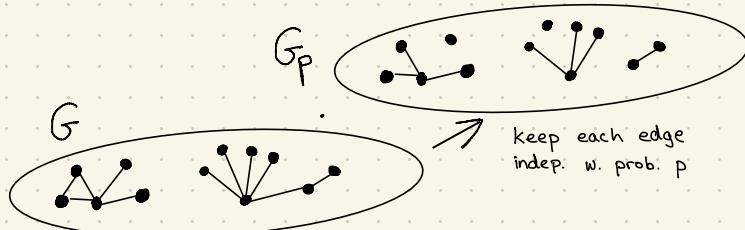
G, G_p similar mod values + partitions

if $\frac{e(G)p}{n} \rightarrow \infty$ whp

$$|q^*(G) - q^*(G_p)| = o(1)$$

$\forall \lambda \exists \lambda':$
(similar)

$$|q_{\lambda'}(G) - q_{\lambda'}(G_p)| = o(1)$$



III Group Exercises

- modularity-based + planted structure
- graph theory, random graphs, theory of algs, simulations