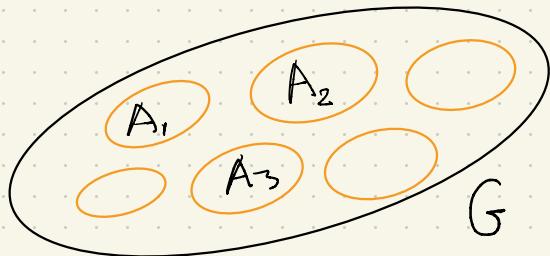


# Modularity + Sampling

with Colin McDiarmid

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



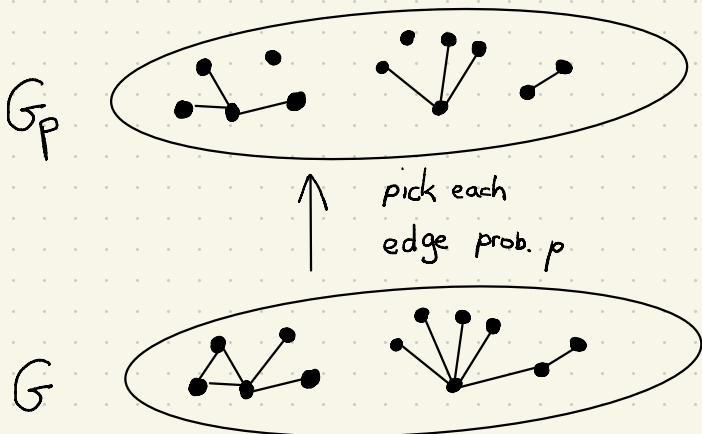
graph  $G$ ,  $m$  edges.  $A = \{A_1, \dots, A_k\}$  vertex partition

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



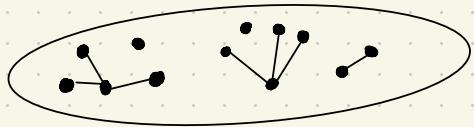
$$G_p = (V, E_p)$$

$E_p$  each edge  
kept indep. prob.  $p$

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

# Sampling

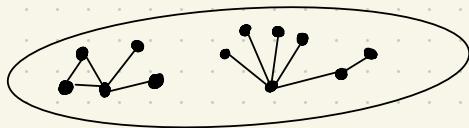
$G_p$



$$G_p = (V, E_p)$$

$E_p$  each edge  
kept indep. prob.  $p$

$G$



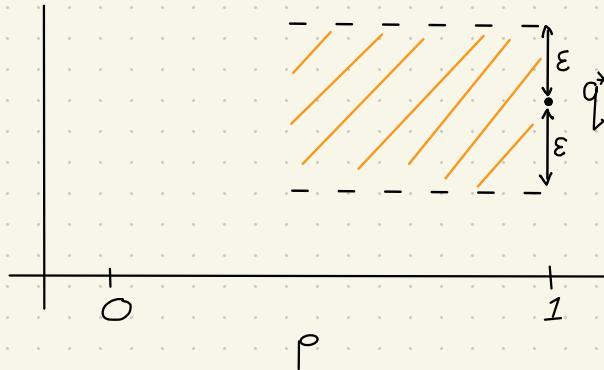
$$G = (V, E) \text{ fixed graph}$$

Q

Given  $G$ ,  $\varepsilon > 0$  for which  $p$  is it true  
with prob  $> 1 - \varepsilon$

$$|q^*(G_p) - q^*(G)| < \varepsilon$$

$q^*(G_p)$

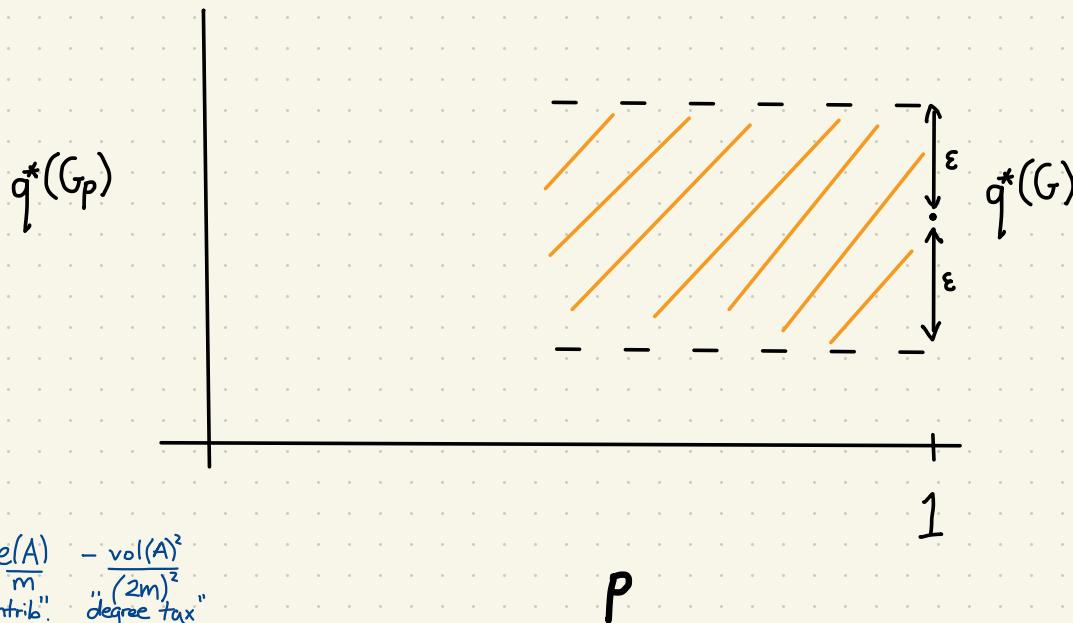


$p$

Thm [McD+S]  $\forall \varepsilon > 0 \exists c = c(\varepsilon)$

for graph  $G$   
constant  $p$

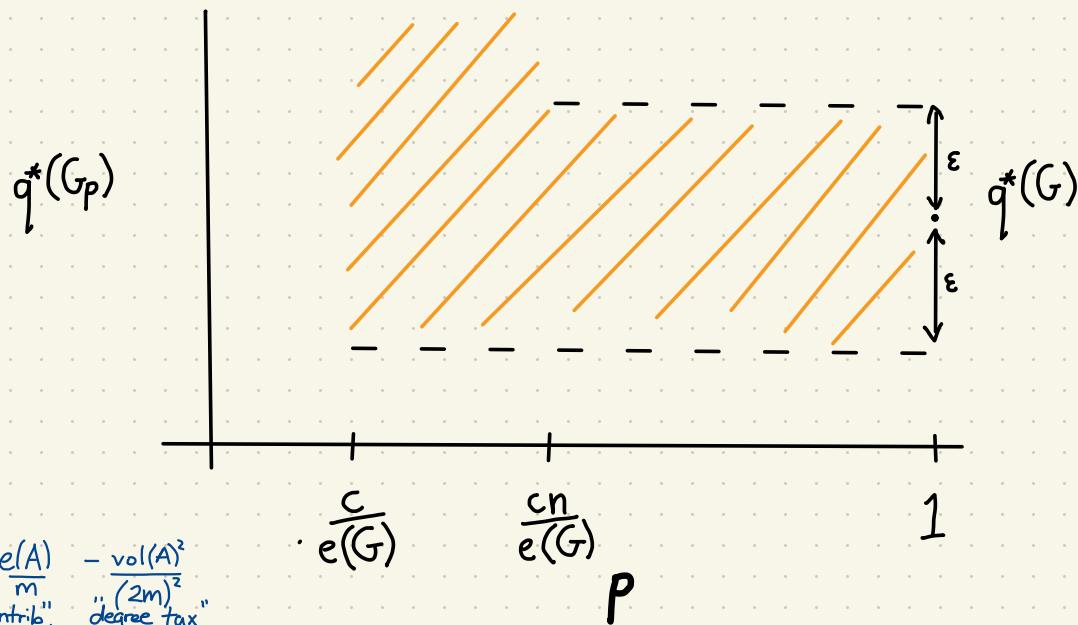
- if                              then w. prob  $> 1 - \varepsilon$        $q^*(G_p) > q^*(G) - \varepsilon$
- if                              "                              "                       $q^*(G_p) < q^*(G) + \varepsilon$



Thm [McD+S]  $\forall \varepsilon > 0 \exists c = c(\varepsilon)$

for graph  $G$   
constant  $p$

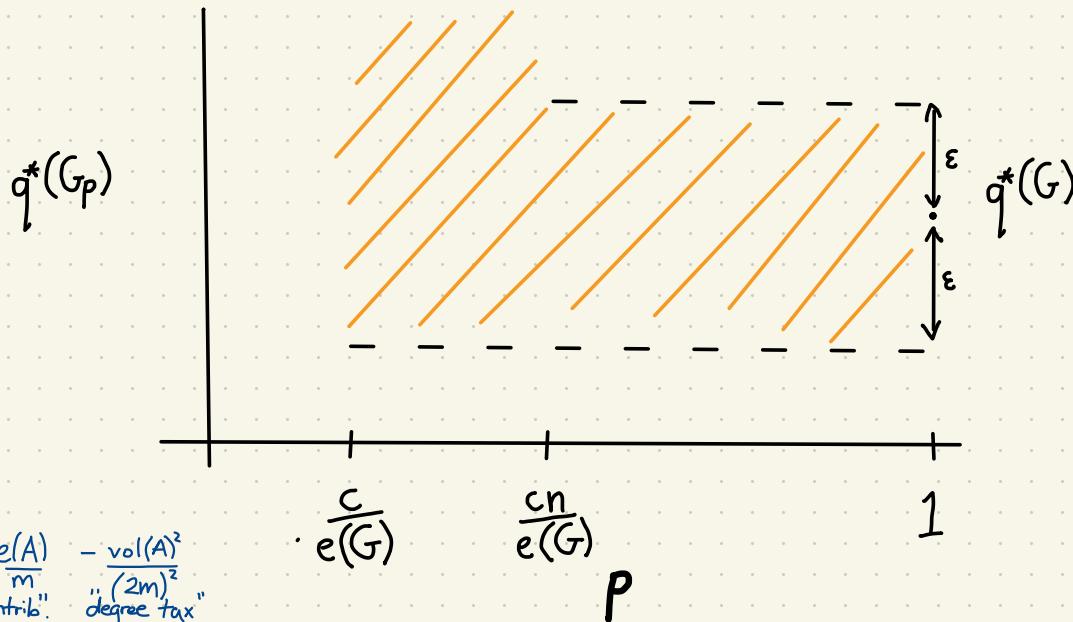
- if  $e(G)p \geq c$  then w. prob  $> 1 - \varepsilon$   $q^*(G_p) > q^*(G) - \varepsilon$
- if  $e(G)p > cn$  " " "  $q^*(G_p) < q^*(G) + \varepsilon$



Thm [McD+S]  $\forall \varepsilon > 0 \exists c = c(\varepsilon)$

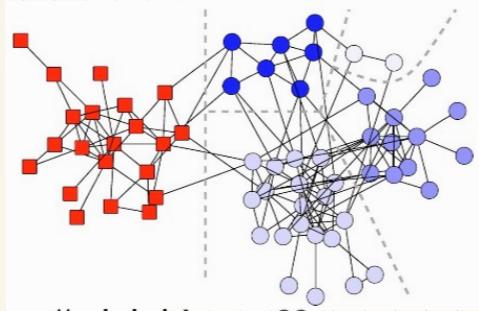
for graph  $G$   
constant  $p$

- if  $e(G)p \geq c$  then w. prob  $> 1 - \varepsilon$   $q^*(G_p) > q^*(G) - \varepsilon$
- if  $e(G)p > cn$  " " "  $q^*(G_p) < q^*(G) + \varepsilon$
- if  $e(G)p \geq cn$  given  $A$ , can construct  $A'$   
w. prob  $> 1 - \varepsilon$   $q_A^*(G) > q_{A'}^*(G_p) - \varepsilon$



## Simulations

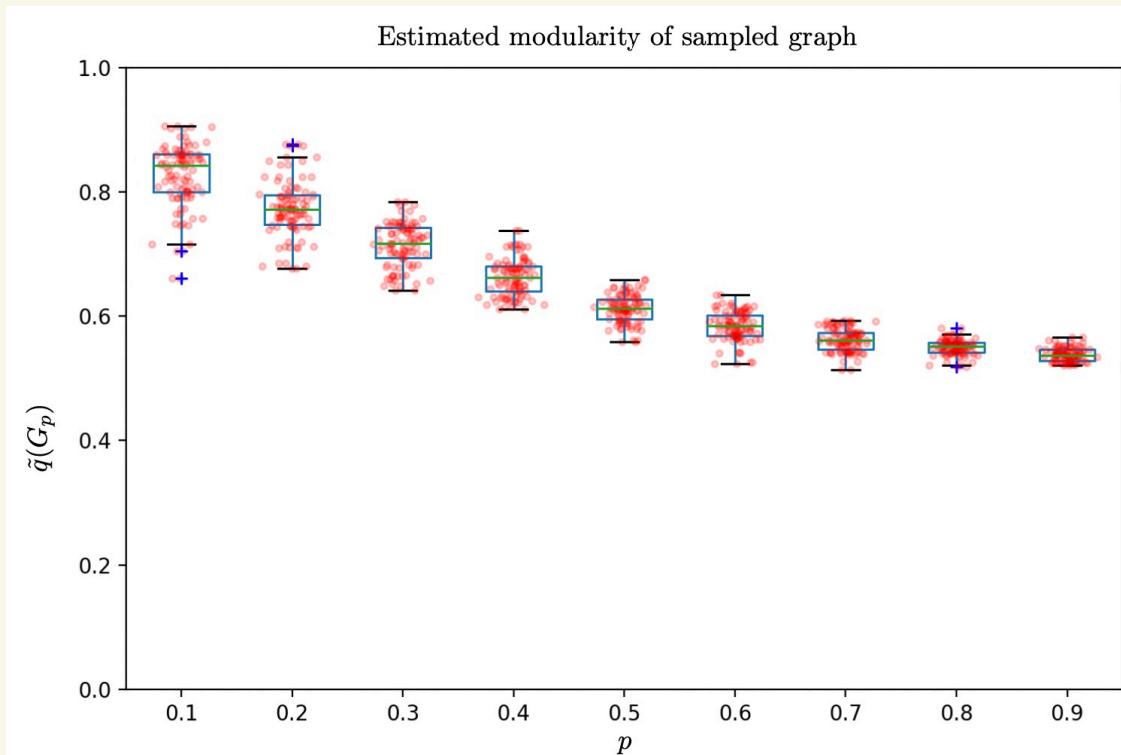
### Dolphin Network [Lusseau]



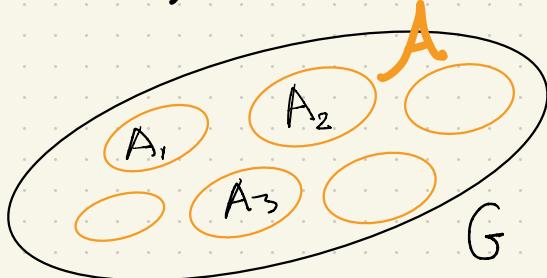
$$|V| = 62 \quad |E| = 152$$

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

[BRANDE + '08]



# 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"

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

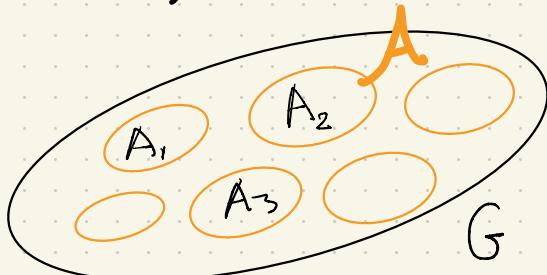
"edge contrib". "degree tax"

$d_u$  = #edges incident to  $u$

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

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

# 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  $G$

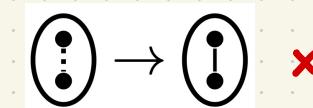
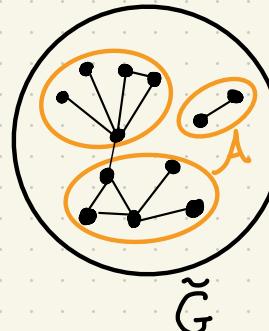
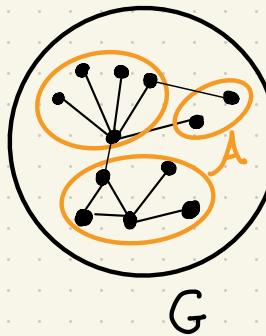
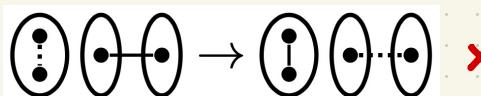
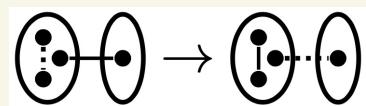
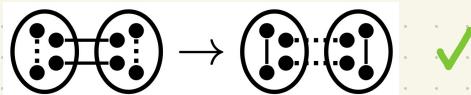
$$q^*(G) = \max_A q_A(G)$$

high vals taken to indicate more community structure"

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

"edge contrib."      "degree tax"

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



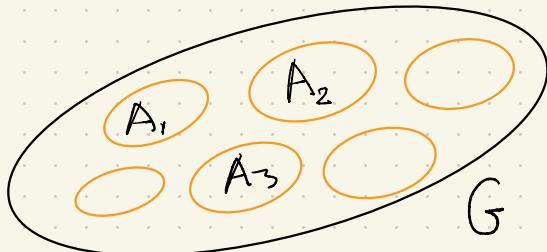
(✓ if  $|A|=2$ ,  $q_A(G) \geq 0$ )

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

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

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

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



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

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

modularity of  $G$

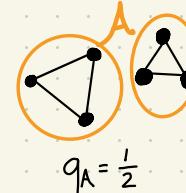
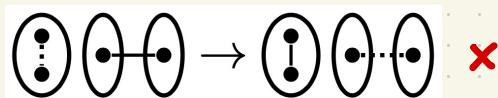
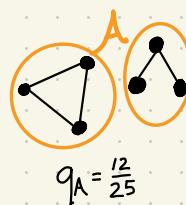
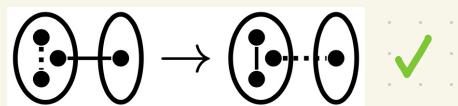
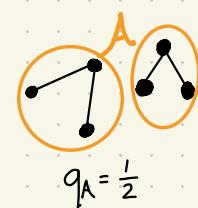
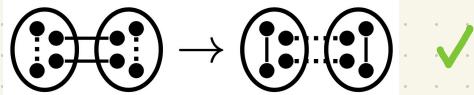
$$q^*(G) = \max_A q_A(G)$$

"high vals taken to indicate more community structure"

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

"edge contrib."      "degree tax"

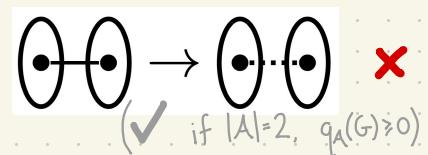
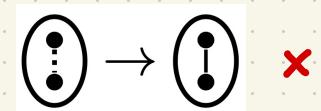
Fix  $A$ , which  $G \rightarrow \tilde{G}$  ensures  $q_A(\tilde{G}) > q_A(G)$  ?



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

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

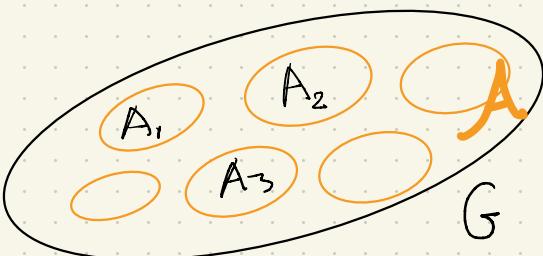
$$\text{vol}(A) = \sum_{u \in A} d_u$$



# Modularity Properties

- Robust to small perturbations in edge set

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



$q^*(G) > 1 - \varepsilon$  if any of following hold

- Connected components in  $G$  all  $< \varepsilon e(G)$  edges

- $\exists A$  with #edges between parts  $< \frac{\varepsilon}{2} e(G)$

and  $\forall A \in A \quad \text{vol}(A) < \frac{\varepsilon}{2} \text{vol}(G)$

- no subgraph  $H \subseteq G$ ,  $H$   $\delta$ -expander,  $e(H) > \delta e(G)$   
(+ Louf)

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib."      "degree tax"

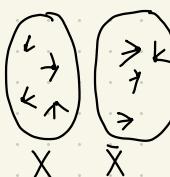
$q^*(G) < \varepsilon$  if any of following hold

- $\bar{\lambda}(G) < \varepsilon$  where  $\bar{\lambda}(G)$  is spectral

gap of Laplacian of  $G$

NB:  $r$ -regular  $G$ ,  $\bar{\lambda}(G) = \frac{1}{r} \max_{i \neq 0} |\lambda_i(A_G)|$

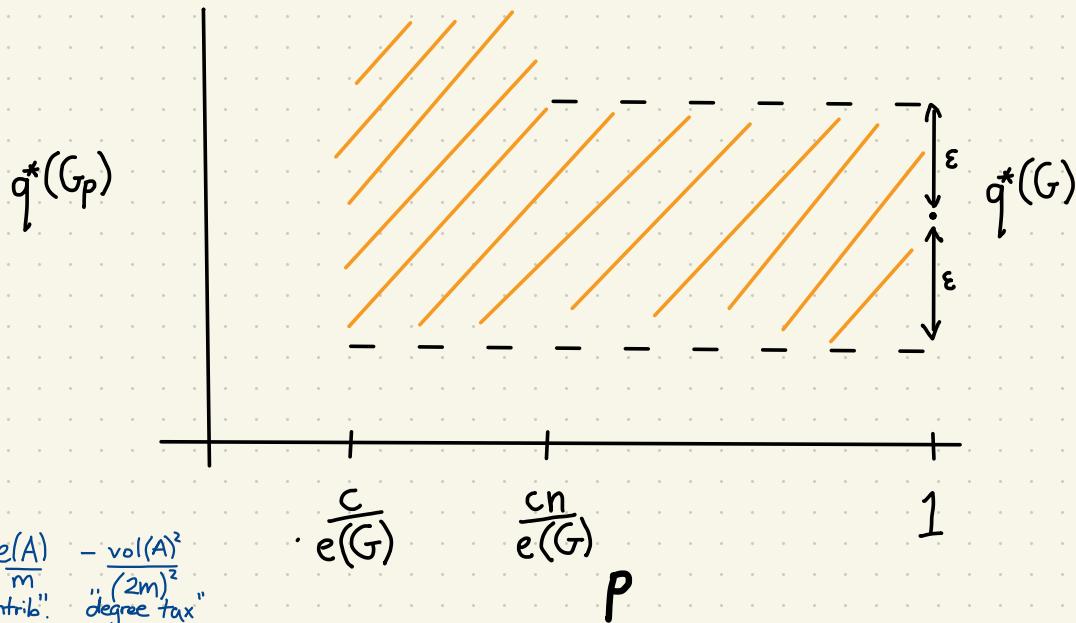
- $\forall X \in V(G), \frac{e(X, \bar{X})}{e(G)} > 2(1-\varepsilon) \frac{\text{vol}(X)}{\text{vol}(G)} \frac{\text{vol}(\bar{X})}{\text{vol}(G)}$



Thm [McD+S]  $\forall \varepsilon > 0 \exists c = c(\varepsilon)$

for graph  $G$   
constant  $p$

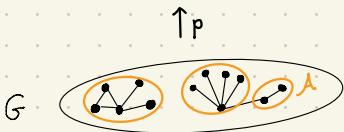
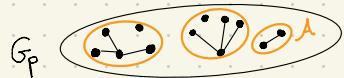
- if  $e(G)p \geq c$  then w. prob  $> 1 - \varepsilon$   $q^*(G_p) > q^*(G) - \varepsilon$
- if  $e(G)p > cn$  " " "  $q^*(G_p) < q^*(G) + \varepsilon$
- if  $e(G)p \geq cn$  given  $A$ , can construct  $A'$   
w. prob  $> 1 - \varepsilon$   $q_A^*(G) > q_{A'}^*(G_p) - \varepsilon$



IDEA OF PROOF: approx opt partition with 'well-behaved' partition

Thm [Dinh+Thai]: for  $k$  positive integer

$$\max_{|A| \leq k} q_A(G) \geq q^*(G)(1 - \frac{1}{k})$$



$A$  is  $\eta$ -fat for  $G$  if  $\text{vol}(A) \geq \eta \text{ vol}(G)$

Fattening Lemma: Given  $A$  partition of  $G$  can construct  $A'$   $\eta$ -fat s.t.  $A'$  refinement of  $A$

$$q_{A'}(G) \geq q_A(G) - 2\eta$$

ASIDE: LOAD BALANCING

$$\vec{x} = (x_1, \dots, x_n) \quad x_i \geq 0 \quad \sum_i x_i = 1$$

$$\text{let } \gamma(\vec{x}) = \max_{I \subseteq [n]} \min \left\{ \sum_{i \in I} x_i, 1 - \sum_{i \in I} x_i \right\}$$

observe

$$\gamma(\vec{x}) \geq \frac{1}{2} - \max_i x_i$$

Q what is the largest  $\alpha$  s.t.

$$\gamma(\vec{x}) \geq \alpha(1 - \sum_i x_i^2)$$

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib." "degree tax"

IDEA OF PROOF: approx opt partition with 'well-behaved' partition

Thm [Dinh+Thai]: for  $k$  positive integer

$$\max_{|A| \leq k} q_A(G) \geq q^*(G)(1 - \frac{1}{k})$$



$\uparrow_P$



$A$  is  $\eta$ -fat for  $G$  if  $\text{vol}(A) \geq \eta \text{ vol}(G)$

Fattening Lemma: Given  $A$  partition of  $G$  can construct  $A'$   $\eta$ -fat s.t.  $A'$  refinement of  $A$

$$q_{A'}(G) \geq q_A(G) - 2\eta$$

$$q_A(G) = \sum_{\substack{A \in A \\ \text{"edge contrib."}}} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

ASIDE: LOAD BALANCING

$$\vec{x} = (x_1, \dots, x_n) \quad x_i \geq 0 \quad \sum_i x_i = 1$$

$$\text{let } \gamma(\vec{x}) = \max_{I \subseteq [n]} \min \left\{ \sum_{i \in I} x_i, 1 - \sum_{i \in I} x_i \right\}$$

observe

$$\gamma(\vec{x}) \geq \frac{1}{2} - \max_i x_i$$

Q what is the largest  $\alpha$  s.t.

$$\gamma(\vec{x}) \geq \alpha(1 - \sum_i x_i^2)$$

$$\alpha \leq \frac{1}{2}$$

$$\gamma\left(\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right) = \frac{1}{3} = \frac{1}{2}\left(1 - \frac{1}{3}\right)$$

$$\alpha \geq \frac{1}{2}$$

greedy alg

- $x_1 \geq \dots \geq x_n$

- init.  $A, B = \emptyset$

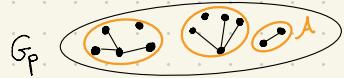
- step i add  $x_i$  to smaller heap

- achieves  $\alpha \geq \frac{1}{2}$

IDEA OF PROOF: approx opt partition with 'well-behaved' partition

Thm [Dinh+Thai]: for  $k$  positive integer

$$\max_{|A| \leq k} q_A(G) \geq q^*(G)(1 - \frac{1}{k})$$



$\uparrow_P$



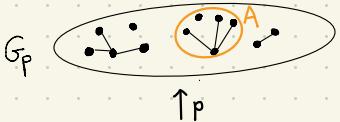
$A$  is  $\eta$ -fat for  $G$  if  $\text{vol}(A) \geq \eta \text{ vol}(G)$

Fattening Lemma: Given  $A$  partition of  $G$  can construct  $A'$   $\eta$ -fat s.t.  $A'$  refinement of  $A$

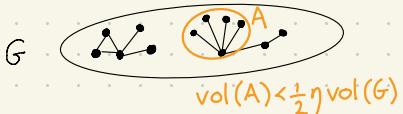
$$q_{A'}(G) \geq q_A(G) - 2\eta$$

## BAD EVENTS

$B_0$ :  $\exists A \subseteq V$  s.t.  $\text{vol}(A) > \eta \text{ vol}(G_P)$

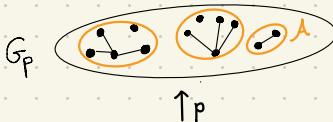


$\uparrow_P$

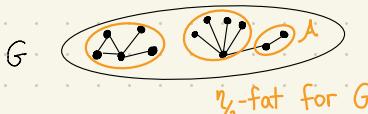


$\text{vol}(A) < \frac{1}{2}\eta \text{ vol}(G)$

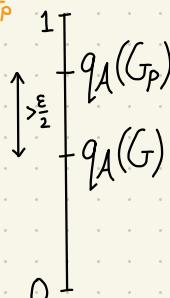
$B_1$ :  $\exists A$  s.t.  $\eta$ -fat for  $G_P$



$\uparrow_P$



$\eta$ -fat for  $G$



$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib."      "degree tax"

## ASIDE: LOAD BALANCING

$$\vec{x} = (x_1, \dots, x_n) \quad x_i \geq 0 \quad \sum_i x_i = 1$$

$$\text{let } \gamma(\vec{x}) = \max_{I \subseteq [n]} \min \left\{ \sum_{i \in I} x_i, 1 - \sum_{i \in I} x_i \right\}$$

observe

$$\gamma(\vec{x}) \geq \frac{1}{2} - \max_i x_i$$

Q what is the largest  $\alpha$  s.t.

$$\gamma(\vec{x}) \geq \alpha(1 - \sum_i x_i^2)$$

$$\alpha \leq \frac{1}{2}$$

$$\gamma\left(\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right) = \frac{1}{3} = \frac{1}{2}\left(1 - \frac{1}{3}\right)$$

$$\alpha \geq \frac{1}{2}$$

greedy alg

- $x_1 \geq \dots \geq x_n$

- init.  $A, B = \emptyset$

- step i add  $x_i$  to smaller heap

achieves  $\alpha \geq \frac{1}{2}$

OPEN - LIMITS: let  $\bar{q}(n, c) = \mathbb{E}(q^*(G_{n, \frac{c}{n}}))$

Conj [MD+S]:  $\forall c > 1$ ,  $\bar{q}(n, c)$  tends to a limit  $\bar{q}(c)$  as  $n \rightarrow \infty$

If conj  
holds then {

- (i) for  $0 < c \leq 1$ ,  $\bar{q}(n, c) \rightarrow \bar{q}(c) = 1$  as  $n \rightarrow \infty$
- (ii)  $0 < \bar{q}(c) < 1$  for  $c > 1$
- (iii)  $\bar{q}(c) = \Theta(c^{-\frac{1}{2}})$  as  $c \rightarrow \infty$
- (iv)  $\bar{q}(c)$  is (uniformly) continuous for  $c \in (0, \infty)$
- (v)  $\bar{q}(c)$  is non-increasing for  $c \in (0, \infty)$

}  
↳ Erdős-Renyi Results  
↳ Sampling Thm

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib." "degree tax"

OPEN - LIMITS: let  $\bar{q}(n, c) = \mathbb{E}(q^*(G_{n, \frac{c}{n}}))$

Conj [MöD+S]:  $\forall c > 1$ ,  $\bar{q}(n, c)$  tends to a limit  $\bar{q}(c)$  as  $n \rightarrow \infty$

- If conj holds then
- (i) for  $0 < c \leq 1$ ,  $\bar{q}(n, c) \rightarrow \bar{q}(c) = 1$  as  $n \rightarrow \infty$
  - (ii)  $0 < \bar{q}(c) < 1$  for  $c > 1$
  - (iii)  $\bar{q}(c) = \Theta(c^{-\frac{1}{2}})$  as  $c \rightarrow \infty$
  - (iv)  $\bar{q}(c)$  is (uniformly) continuous for  $c \in (0, \infty)$
  - (v)  $\bar{q}(c)$  is non-increasing for  $c \in (0, \infty)$

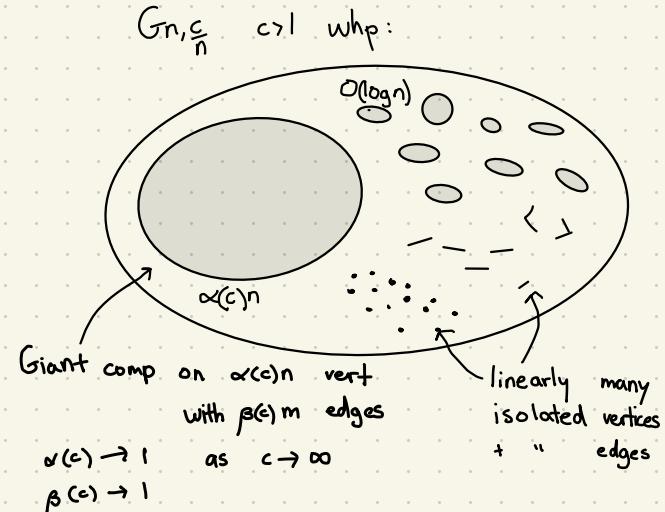
$\left. \begin{array}{l} \text{Erdős-Renyi Results} \\ \text{Sampling Thm} \end{array} \right\}$

## OPEN - FIVE PARTS

Conj [Reichardt+Bornholdt '06]:  $\max_{|A|=5} q_A(G_{n, \frac{c}{n}}) = q^*(G_{n, \frac{c}{n}})$  whp

$$q_A(G) = \sum_{A \in A} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib." "degree tax"



OPEN - LIMITS: let  $\bar{q}(n, c) = \mathbb{E}(q^*(G_{n, \frac{c}{n}}))$

Conj [M'D+S]:  $\forall c > 1$ ,  $\bar{q}(n, c)$  tends to a limit  $\bar{q}(c)$  as  $n \rightarrow \infty$

If conj  
holds then

- (i) for  $0 < c \leq 1$ ,  $\bar{q}(n, c) \rightarrow \bar{q}(c) = 1$  as  $n \rightarrow \infty$
- (ii)  $0 < \bar{q}(c) < 1$  for  $c > 1$
- (iii)  $\bar{q}(c) = \Theta(c^{-\frac{1}{2}})$  as  $c \rightarrow \infty$
- (iv)  $\bar{q}(c)$  is (uniformly) continuous for  $c \in (0, \infty)$
- (v)  $\bar{q}(c)$  is non-increasing for  $c \in (0, \infty)$

$\left. \begin{array}{l} \text{Erdős-Renyi Results} \\ \text{Sampling Thm} \end{array} \right\}$

## OPEN - FIVE PARTS

Conj [Reichardt+Bornholdt '06]:  $\max_{|A|=5} q_A(G_{n, \frac{c}{n}}) = q^*(G_{n, \frac{c}{n}})$  whp

But!  $\forall c \exists \delta(c)$   $\max_{|A|=k} q_A(G_{n, \frac{c}{n}}) \leq q^*(G_{n, \frac{c}{n}})(1 - \delta(c))$  whp

By Dinh+Thai  $\max_{|A| \leq k} q_A(G_{n, \frac{c}{n}}) \geq q^*(G_{n, \frac{c}{n}})(1 - \frac{1}{k})$

Conj [M'D+S]:  $\exists K$  ( $k=5?$ ) st.  $\forall \varepsilon > 0 \exists c_0$  and  $\forall c > c_0$

$$\max_{|A| \leq k} q_A(G_{n, \frac{c}{n}}) \geq q^*(G_{n, \frac{c}{n}})(1 - \varepsilon)$$

$$q_A(G) = \sum_{A \subseteq E} \frac{e(A)}{m} - \frac{\text{vol}(A)^2}{(2m)^2}$$

"edge contrib." "degree tax"

