

# Goodness of Fit Tests for Random Multigraph Models

Termeh Shafie  
Department of Social Statistics  
The Mitchell Centre for Social Network Analysis  
The University of Manchester

*graphs where **multiple edges** and **self-edges** are permitted*

- can appear directly in applications (although scarce)
- can be constructed by different kinds of aggregations in graphs
  - aggregation based on vertex attributes
  - aggregation based on edge attributes



- multigraphs represented by their edge multiplicity sequence

$$\mathbf{m} = (m_{ij} : (i, j) \in R)$$

where  $R$  is the canonical site space for undirected edges

$$R = \{(i, j) : 1 \leq i \leq j \leq n\}$$

$$(1, 1) < (1, 2) < \dots < (1, n) < (2, 2) < (2, 3) < \dots < (n, n)$$

- the number of vertex pair sites is given by  $r = \binom{n+1}{2}$
- edge multiplicities as entries in a matrix

$$\mathbf{m} = \begin{bmatrix} m_{11} & m_{12} & \dots & m_{1n} \\ 0 & m_{22} & \dots & m_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & m_{nn} \end{bmatrix} \quad \mathbf{m} + \mathbf{m}' = \begin{bmatrix} 2m_{11} & m_{12} & \dots & m_{1n} \\ m_{12} & 2m_{22} & \dots & m_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ m_{1n} & m_{2n} & \dots & 2m_{nn} \end{bmatrix}$$

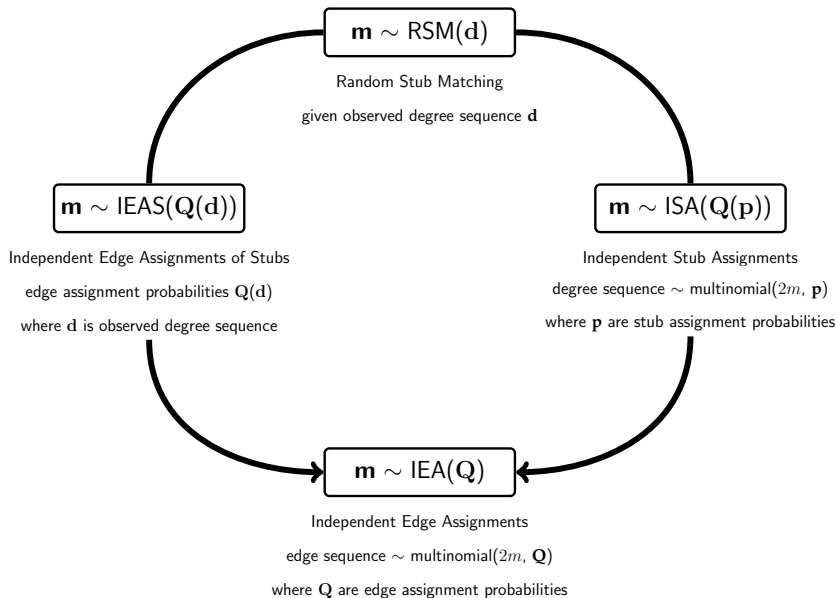
## 1. random stub matching (RSM)

- edges are assigned to sites given fixed degree sequence  $\mathbf{d} = (d_1, \dots, d_n)$
- probability that an edge is assigned to site  $(i, j) \in R$

$$Q_{ij} = \begin{cases} \binom{d_i}{2} / \binom{2m}{2} & \text{for } i = j \\ d_i d_j / \binom{2m}{2} & \text{for } i < j \end{cases}$$

## 2. independent edge assignment (IEA)

- edges are independently assigned to vertex pairs in site space  $R$
- edge assignment probabilities  $\mathbf{Q} = (Q_{ij} : (i, j) \in R)$
- $\mathbf{m}$  is multinomial distributed with parameters  $m$  and  $\mathbf{Q}$
- statistics for analysing local and global structure are easily derived
- two variants:
  - **independent edge assignment of stubs (IEAS)**
  - **independent stub assignment (ISA)**



gof measures between observed and expected edge multiplicity sequence  
under simple or composite hypothesis

- test statistics:  $S$  of Pearson and  $A$  of information divergence type
- expected values of the Pearson statistic are derived
- exact distributions of the test statistics are numerically investigated

answers sought to the following:

- are significance levels of test statistics for small number of edges far from those of the asymptotic distribution?
- is the convergence of the cdf's of test statistics slow or rapid?
- does the convergence speed depend on specific parameters in models?
- can better approximations to the actual distributions be obtained using adjustments of the  $\chi^2$ -distributions?
- can power approximations be made for small number of edges?
- how does RSM influence the distributions of statistics?
- how can RSM be tested?

edge multiplicities according to IEA( $\mathbf{Q}$ ) and correct model  $\mathbf{Q}_0 = \mathbf{Q}$  tested:

- the Pearson statistic

$$S_0 = \sum_{i \leq j} \sum \frac{(m_{ij} - mQ_{0ij})^2}{mQ_{0ij}} = \sum_{i \leq j} \sum \frac{m_{ij}^2}{mQ_{0ij}} - m \stackrel{asympt}{\sim} \chi^2(r-1)$$

- the divergence statistic

$$D_0 = \sum_{i \leq j} \sum \frac{m_{ij}}{m} \log \frac{m_{ij}}{mQ_{0ij}} \quad \text{and} \quad A_0 = \frac{2m}{\log e} D_0 \stackrel{asympt}{\sim} \chi^2(r-1)$$

the composite multigraph hypothesis

- ISA for unknown  $\mathbf{p}$
- IEAS for unknown  $\mathbf{d}$

parameters have to be estimated from data  $\mathbf{m}$

when correct model is tested:

- the Pearson statistic

$$\hat{S} = \sum_{i \leq j} \sum \frac{(m_{ij} - m\hat{Q}_{ij})^2}{m\hat{Q}_{ij}} = \sum_{i \leq j} \sum \frac{m_{ij}^2}{m\hat{Q}_{ij}} - m \stackrel{asympt}{\sim} \chi^2(r - n)$$

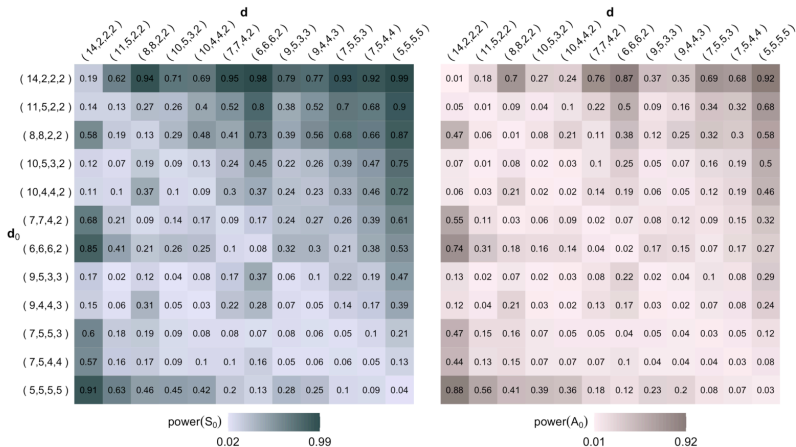
- the divergence statistic

$$\hat{D} = \sum_{i \leq j} \sum \frac{m_{ij}}{m} \log \frac{m_{ij}}{m\hat{Q}_{ij}} \quad \text{and} \quad \hat{A} = \frac{2m}{\log e} \hat{D} \stackrel{asympt}{\sim} \chi^2(r - n)$$

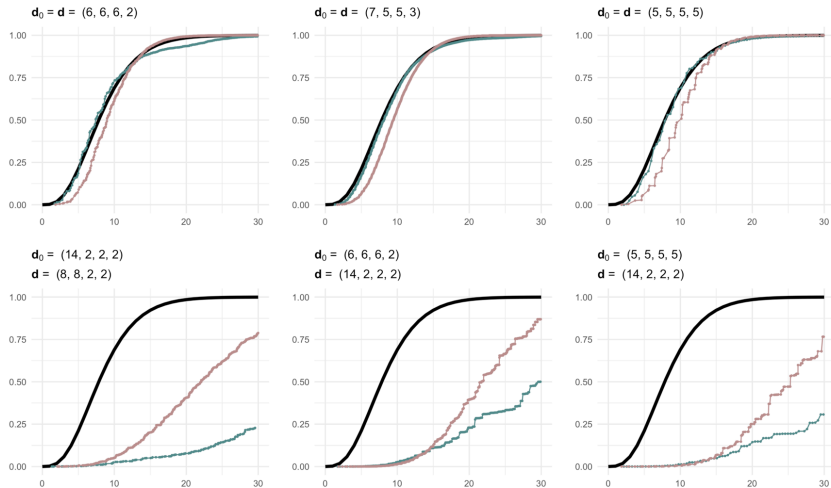


# goodness of fit: test illustrations

power when simple IEAS( $\mathbf{d}_0$ ) hypotheses are tested against IEAS( $\mathbf{d}$ ) models for multigraphs with  $n = 4$ ,  $m = 10$  and  $\alpha(\chi^2_9) = 0.04$

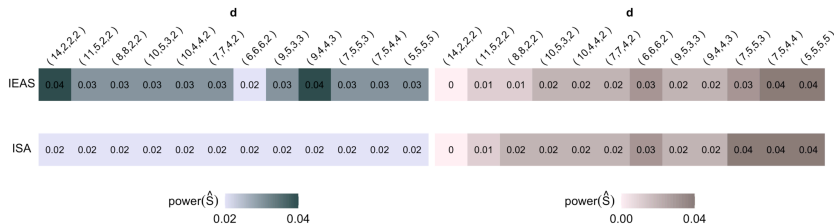


null and non-null distributions of  $S_0$  and  $A_0$ , and the  $\chi^2_9$ -distribution when simple IEAS( $\mathbf{d}_0$ ) hypotheses are tested against IEAS( $\mathbf{d}$ )

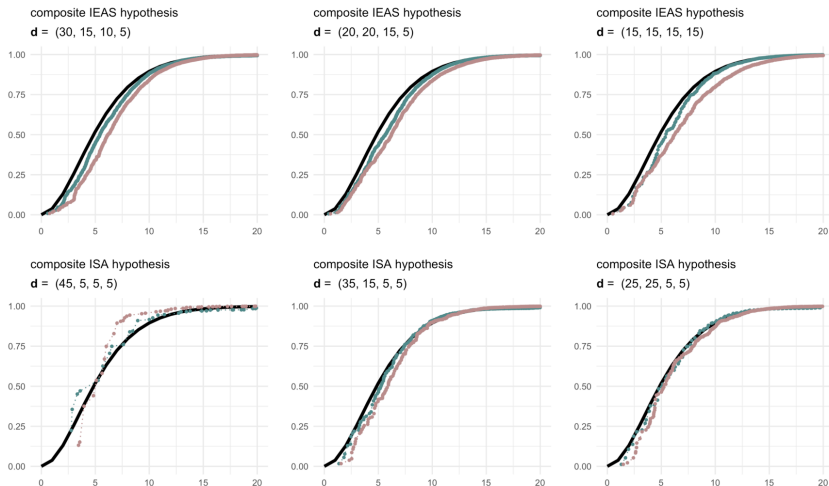


# goodness of fit: test illustrations

probabilities of false rejection (top) and power (bottom) when composite IEAS and ISA hypotheses are tested against IEAS(**d**) models for multigraphs with  $n = 4$ ,  $m = 10$  and  $\alpha(\chi_6^2) = 0.04$



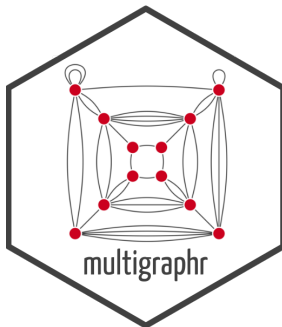
non-null distributions of  $\hat{S}$ ,  $\hat{A}$ , and the  $\chi^2_6$ -distribution when composite IEAS and ISA hypotheses tested against RSM(**d**) models



# summary of error probabilities

model		simple IEAS( $d_0$ ) hypothesis			composite hypothesis	
		$d_0 = d$	Flat $d_0 \neq d$	Skew $d_0 \neq d$	IEAS	ISA
IEAS	Flat $d$	$\alpha_{S_0} > \alpha_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\alpha_{\hat{S}} \leq \alpha_{\hat{A}}$	$\beta_{\hat{S}} > \beta_{\hat{A}}$
	Skew $d$	$\alpha_{S_0} > \alpha_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\alpha_{\hat{S}} > \alpha_{\hat{A}}$	$\beta_{\hat{S}} \geq \beta_{\hat{A}}$
		simple ISA( $d_0/2m$ ) hypothesis			composite hypothesis	
		$d_0 = d$	Flat $d_0 \neq d$	Skew $d_0 \neq d$	IEAS	ISA
ISA	Flat $d$	$\alpha_{S_0} \geq \alpha_{A_0}$	$\beta_{S_0} \leq \beta_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	inconclusive	$\alpha_{\hat{S}} \leq \alpha_{\hat{A}}$
	Skew $d$	$\alpha_{S_0} > \alpha_{A_0}$	$\beta_{S_0} \leq \beta_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\beta_{\hat{S}} < \beta_{\hat{A}}$	$\alpha_{\hat{S}} > \alpha_{\hat{A}}$
		simple IEAS( $d_0$ ) or ISA( $d_0/2m$ ) hypothesis			composite hypothesis	
		$d_0 = d$	Flat $d_0 \neq d$	Skew $d_0 \neq d$	IEAS	ISA
RSM	Flat $d$	$\beta_{S_0} \geq \beta_{A_0}$	inconclusive	$\beta_{S_0} = \beta_{A_0}$	$\beta_{\hat{S}} > \beta_{\hat{A}}$	$\beta_{\hat{S}} > \beta_{\hat{A}}$
	Skew $d$	$\beta_{S_0} \leq \beta_{A_0}$	$\beta_{S_0} = \beta_{A_0}$	$\beta_{S_0} < \beta_{A_0}$	$\beta_{\hat{S}} \geq \beta_{\hat{A}}$	$\beta_{\hat{S}} > \beta_{\hat{A}}$

- even for very small  $m$ , the null distributions of the test statistics under IEA are fairly well approximated by their asymptotic distributions
- the convergence of the cdf's of test statistics are rapid and depend on the parameters in models
- approximations to the actual distributions can be obtained using adjustments of the  $\chi^2$ -distributions yielding better power
- the influence of RSM on both test statistics is substantial for small  $m$ , implying a shift of their distributions towards smaller values compared to what holds true for the null distributions under IEA



<https://github.com/termehs/multigraphr>

```
# install.packages("devtools")  
devtools::install_github("termehs/multigraphr")
```