

# Estimating Subgraph Generation Models to Understand Large Network Formation

1<sup>st</sup> Laurens Bogaardt  
Netherlands eScience Center  
Amsterdam, the Netherlands  
l.bogaardt@esciencecenter.nl

2<sup>nd</sup> Frank W. Takes  
University of Amsterdam  
Amsterdam, the Netherlands  
takes@uva.nl

**Abstract**—Recently, a new network formation model was proposed: SUGM. Our research looks into a method to estimate the parameters of this model based on the subgraph census.

**Index Terms**—Networks, Graphs, ERGM, SUGM, Subgraphs

Social scientists often aim to understand the incentives and mechanisms which result in large scale socio-economic structures. Key to this is network formation analysis. However, large datasets are not uncommon, leading to a computational challenge. For example, political scientists interested in global networks of corporate control may need to analyse millions of companies to answer their questions [1], [2].

The Exponential Random Graph Model (ERGM) is a frequently used network formation model. Unfortunately, it suffers from two fundamental flaws. Firstly, its parameter estimates are inconsistent [3], [4]. Secondly, it does not scale well [5]. Recently, an alternative network formation model was suggested: the Subgraph Generation Model (SUGM) [6]–[8].

A SUGM is defined by a set of  $l$  small subgraphs, such as links, triangles or stars, each with corresponding probabilities. For each subgraph  $i$  of  $m_i$  nodes, the  $n$  nodes of the entire network are grouped into all possible subsets of  $m_i$  nodes. Then, each of these subsets receives the subgraph  $i$  with probability  $1 - p_i$  or remains empty with probability  $p_i$ .

The observed network, left in Fig. 1, is the union of all these subgraphs, right in Fig. 1, where the generated subgraphs may overlap. Multiple neighbouring subgraphs may incidentally form additional structures such as triangles or squares.

The original articles describing SUGM contain two methods

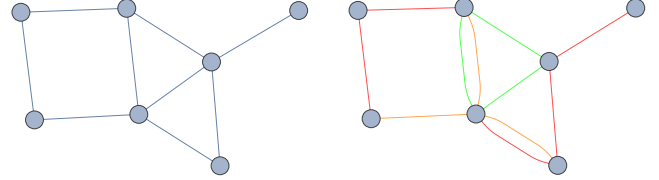


Fig. 1. The observed network (left) is the union (right) of randomly generated links (red), triangles (green) and stars (yellow).

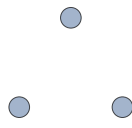
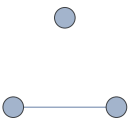
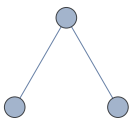
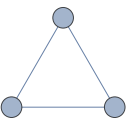
to estimate the parameters of the model. Our work suggests a third, more intuitive method based on the subgraph census. In a  $k$ -subgraph census, a network of  $n$  nodes is grouped into all possible subsets of  $k$  nodes, which are then tallied according to their isomorphism class [9]–[11]. This counting can be achieved in polynomial time.

Table I contains the probabilities of observing any of the possible triads for three different generation models. In general, each of the  $r$  counts of the census  $x_j$ , together with the probability functions  $f_j(\hat{p}_1, \dots, \hat{p}_l)$ , enter into the multinomial probability mass function of (1) to form the likelihood function. This can be used to estimate the parameters of the model and their confidence intervals.

$$\mathcal{L}(f_1, \dots, f_r | x_1, \dots, x_r) = \frac{\Gamma(\sum_j x_j + 1)}{\prod_j \Gamma(x_j + 1)} \prod_{j=1}^r f_j^{x_j} \quad (1)$$

Future work should extend the list of possible subgraphs, deal with the correlations within the census, develop an  $R$ -package and apply the model to real-world data.

TABLE I  
PROBABILITIES IN THE SUBGRAPH CENSUS

Model	Subgraphs of the Undirected Triad Census			
				
Links	$p_L^3$	$3 p_L^2 (1 - p_L)$	$3 p_L (1 - p_L)^2$	$(1 - p_L)^3$
Triangles	$p_T (p_T^{n-3})^3$	$3 p_T (p_T^{n-3})^2 (1 - p_T^{n-3})$	$3 p_T (p_T^{n-3}) (1 - p_T^{n-3})^2$	$(1 - p_T) + p_T (1 - p_T^{n-3})^3$
Links & Triangles	$p_T (p_L p_T^{n-3})^3$	$3 p_T (p_L p_T^{n-3})^2 (1 - p_L p_T^{n-3})$	$3 p_T (p_L p_T^{n-3}) (1 - p_L p_T^{n-3})^2$	$(1 - p_T) + p_T (1 - p_L p_T^{n-3})^3$

## REFERENCES

- [1] F. W. Takes and E. M. Heemskerk, "Centrality in the global network of corporate control," *Social Network Analysis and Mining*, vol. 6, pp. 1–18, 2016.
- [2] M. Fennema and E. M. Heemskerk, "When theory meets methods: the naissance of computer assisted corporate interlock research," *Global Networks*, vol. 1, pp. 81–104, 2018.
- [3] C. R. Shalizi and A. Rinaldo, "Consistency under sampling of exponential random graph models," *The Annals of Statistics*, vol. 41, pp. 508–535, 2013.
- [4] S. Chatterjee and P. Diaconis, "Estimating and understanding exponential random graph models," *The Annals of Statistics*, vol. 41, pp. 2428–2461, 2013.
- [5] S. Bhamidi, G. Bresler, and A. Sly, "Mixing time of exponential random graphs," *The Annals of Applied Probability*, vol. 21, pp. 2146–2170, 2011.
- [6] A. G. Chandrasekhar and M. O. Jackson, "Tractable and consistent random graph models," *ArXiv*, 2014. [Online]. Available: <https://arxiv.org/abs/1210.7375>
- [7] —, "A network formation model based on subgraphs," *ArXiv*, 2016. [Online]. Available: <https://arxiv.org/abs/1611.07658>
- [8] A. G. Chandrasekhar, "Econometrics of network formation," in *The Oxford Handbook of the Economics of Networks*, Y. Bramoulle, A. Galeotti, and B. Rogers, Eds. Oxford University Press, 2016.
- [9] J. A. Davis and S. Leinhardt, "The structure of positive interpersonal relations in small groups," in *Sociological Theory in Progress*, J. Berger, M. Zelditch, and B. Anderson, Eds. Houghton-Mifflin, 1972.
- [10] P. W. Holland and S. Leinhardt, "A method for detecting structure in sociometric data," *American Journal of Sociology*, vol. 76, pp. 492–513, 1970.
- [11] —, "Local structure in social networks," *Sociological Methodology*, vol. 7, pp. 1–45, 1976.