

Generalised hypergeometric ensembles of random graphs: the configuration model as an urn problem

Giona Casiraghi* Vahan Nanumyan[†]

*Chair of Systems Design,
ETH Zurich, Weinbergstrasse 56/58, 8092 Zurich, Switzerland*

**gcasiraghi@ethz.ch*

[†]*vnanumyan@ethz.ch*

Abstract

We introduce a broad class of random graph models: the generalised hypergeometric ensemble (GHypEG). This class enables to solve some long standing problems in random graph theory. First, GHypEG provides an elegant and compact formulation of the well-known configuration model in terms of an urn problem. Second, GHypEG allows to incorporate arbitrary tendencies to connect different vertex pairs. Third, we present the closed-form expressions of the associated probability distribution ensures the analytical tractability of our formulation. This is in stark contrast with the previous state-of-the-art, which is to implement the configuration model by means of computationally expensive procedures.

Keywords: random graph, network ensemble, configuration model, Wallenius' non-central hypergeometric distribution

1 Introduction

Important features of real-world graphs are often analysed by studying the deviations of empirical observations from suitable random models. Such models, or graph ensembles, are built so that some of the properties of the analysed empirical graph are preserved. Then, one can identify which other features of the empirical graph can be expected at random and which cannot, given the encoded constraints.

The simplest random graph model, named after Erdős and Rényi, generates edges between a given number of vertices with a fixed probability $p \in (0,1]$ [10]. In this model, the properties of vertices, such as their degrees, have all the same expected

value. However, most empirical graphs have heterogeneous, heavy tailed degree distributions [1, 13, 8, 2]. Hence, random graph models able to incorporate arbitrary degree sequences are of special importance.

The most common model fixing only degree sequences is known as the *configuration model* of random graphs [8, 9, 3, 14]. The comparison of an empirical graph with the corresponding configuration model then allows to quantify which properties of the original graph can be ascribed to the degree sequence. The properties not explained by the degree sequence highlight the unique structure of the studied empirical graph.

The standard configuration model has a crucial drawback: the lack of analytical tractability. In fact, the model is realised by means of a repeated rewiring procedure. Each vertex is assigned a number of half-edges, or stubs, corresponding to their degree and one random realisation of the model is obtained by wiring pairs of stubs together uniformly at random. This is a computationally expensive procedure, which does not allow to explore the whole probability space of the model. This problem is exacerbated in the case of larger graphs and graphs with highly heterogeneous degree distributions.

Moreover, the standard configuration model is limited to the coarse analysis of the combinatorial randomness arising from vertex degrees, as it has no free parameters. While it has been invaluable for the macroscopic and the mesoscopic analysis of graphs, such as for graph partitioning through modularity maximisation [16] and for quantifying degree correlations [17], other graph models are needed to address complex dyadic patterns beyond degrees.

In this article, we propose a novel analytically tractable model for random graphs with given expected degree sequences. The formulation of the model relies on mapping the process of drawing edges to a multivariate urn problem. In its simplest case, our model corresponds to the configuration model for directed or undirected multi-edge graphs, fixing in expectation the values of vertex degrees instead of their exact values. In the general case, the model incorporates a parameter for each pair of vertices, which we call *edge propensity*. This parameter controls the relative likelihood of drawing an individual edge between the respective pair of vertices, as opposed to any other pair. This is achieved by biasing the combinatorial edge drawing process. The proposed formulation allows to model and test for arbitrary graph patterns that can be reduced to dyadic relations.

2 Results

Let us consider a *multi-graph* $\mathcal{G} = (V, E)$, where V is a set of n vertices, and $E \subseteq V \times V$ is a multi-set of m (directed or undirected) multi-edges. For clarity, we first provide working definitions of multi-edges and multi-graphs.

Definition 1 (Multi-edges). Let $V \subset \mathbb{N}$ be a set of vertices and $i, j \in V$ two vertices. Elements $(i, j)_l \in E$ and $(i, j)_k \in E$, $l \neq k$, incident to the same two vertices are called multi-edges. The number $A_{ij} \in \mathbb{N}_0$ of multi-edges incident to the same two vertices i and j defines the *multiplicity* of the edge (i, j) .

Definition 2 (Multi-graph). Let $V \subset \mathbb{N}$ be a set of vertices. A graph $\mathcal{G}(V, E)$ is called a multi-graph if $E \subseteq V \times V$ is a *multi-set* of $m := |E|$ multi-edges. Self-loops $(i, i) \in E$ for $i \in V$ are generally allowed. A multi-graph \mathcal{G} can be directed or undirected.

We indicate with \mathbf{A} the adjacency matrix of the graph where entries $A_{ij} \in \mathbb{N}_0$ capture the multiplicity of an edge $(i, j) \in V \times V$ in the multi-set E . In the case of undirected graphs, the adjacency matrix is symmetric, i.e. $\mathbf{A} = \mathbf{A}^T$, and the elements on its diagonal equal twice the multiplicity of the corresponding self-loops.

Definition 3 (Degrees). For each vertex $i \in V$ the in-degree $k_i^{\text{in}} := \sum_{j \in V} A_{ji}$ and the out-degree $k_i^{\text{out}} := \sum_{j \in V} A_{ij}$. The total number of multi-edges is expressed as $m = \sum_{i, j \in V} A_{ij} = \sum_{i \in V} k_i^{\text{out}} = \sum_{j \in V} k_j^{\text{in}}$. We denote the in-degree and out-degree *sequence* of a directed graph \mathcal{G} as $\mathbf{k}^{\text{in}}(\mathcal{G}) = \{k_i^{\text{in}}\}_{i \in V}$ and $\mathbf{k}^{\text{out}}(\mathcal{G}) = \{k_i^{\text{out}}\}_{i \in V}$. We denote with $k_i^{\text{in/out}}(\mathcal{G})$ the i -th entry of the degree sequence $\mathbf{k}^{\text{in/out}}(\mathcal{G})$, corresponding to the in- or out-degree of vertex i . For undirected graphs, the adjacency matrix is symmetric and thus $k_i^{\text{in}} = k_i^{\text{out}} =: k_i$. Hence, there is one degree sequence of an undirected graph, which we denote $\mathbf{k}(\mathcal{G})$.

As we only deal with multi-graphs, we will refer to multi-graphs simply as graphs in the rest of the article.

2.1 Soft Configuration Model

The concept underlying our random graph model is the same as for the standard configuration model of Molloy and Reed [14, 15], which is to randomly shuffle the multi-edges of a graph \mathcal{G} while preserving vertex degrees. The standard configuration model generates multi-edges one after another by sampling uniformly at random a vertex

with an available out-stub (outwards half-edge) and a vertex with an available in-stub (inwards half-edge), until all stubs are consumed. Figure 1 illustrates one step of this process. The resulting random graphs all have exactly the same degree sequence as the original graph \mathcal{G} . All pairs of available in- and out-stubs are equiprobable to be picked, so are the corresponding individual multi-edges. Therefore, the probability to observe a multi-edge between a given pair of vertices positively relates to the number of possible stub pairings of the two vertices, which in turn is defined by the corresponding degrees of these. In fact, this probability depends only on the degrees of the two vertices and on the total number of multi-edges in the graph.

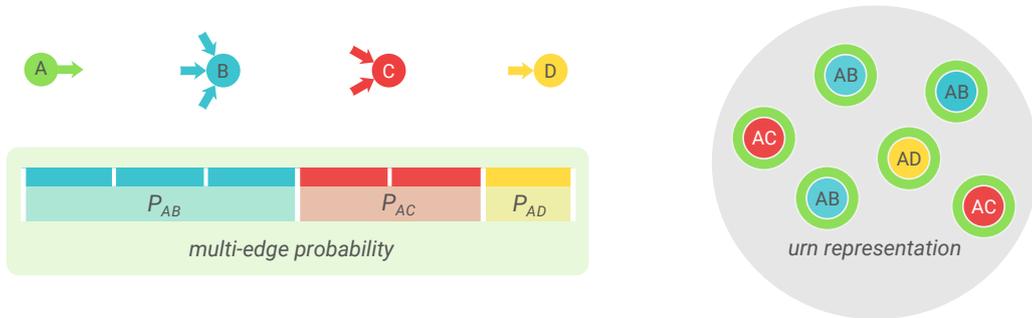


Figure 1: The configuration model represented (**upper left**) as a conventional edge rewiring process and (**right**) as an urn problem. In the former case, once the out-stub (A, \cdot) has been sampled for rewiring, then one in-stub is sampled uniformly at random from those available, to draw a new multi-edge. If we represent each possible combination of an out-stub and an in-stub as a ball, we arrive at the urn problem without replacement. For the shown vertices, the odds of observing a multi-edge (A, B) are three times higher than of observing a multi-edge between (A, D) and 1.5 times higher than of observing a multi-edge between (A, C) in both model representations.

The need to consequently sample two vertices at each step makes it cumbersome to analytically formulate the procedure described above. To overcome this challenge, we take an *edge-centric approach* of sampling m multi-edges from a certain larger multi-set, which we define below in Definition 4. As a consequence of this change of perspective the model will preserve the *expected* degree sequences instead of the exact ones. To this end, we introduce the definition of the *soft configuration model*.

For each pair of vertices $i, j \in V$, we define the number Ξ_{ij} of stub combinations that exist between vertices i and j , which can be conveniently represented in matrix form:

Definition 4 (Combinatorial matrix). We define combinatorial matrix $\Xi \in \mathbb{N}^n \times \mathbb{N}^n$ for graph \mathcal{G} as the matrix whose elements Ξ_{ij} are defined as

$$\Xi_{ij} = k_i^{\text{out}}(\mathcal{G})k_j^{\text{in}}(\mathcal{G}) \quad \text{for } i, j \in V, \quad (1)$$

where $\mathbf{k}^{\text{out}}(\mathcal{G})$ and $\mathbf{k}^{\text{in}}(\mathcal{G})$ are the out-degree and in-degree sequences of the graph \mathcal{G} .

2.1.1 Soft configuration model for directed graphs

Definition 5 (Directed soft configuration model). Let $\hat{\mathbf{k}}^{\text{in}}, \hat{\mathbf{k}}^{\text{out}} \in \mathbb{N}^n$ be in- and out-degree sequences and \hat{V} a set of n vertices. The soft configuration model X generated by $(\hat{V}, \hat{\mathbf{k}}^{\text{in}}, \hat{\mathbf{k}}^{\text{out}})$ is the n^2 -dimensional random vector X defined on the probability space (S, P) with sample space

$$S = \left\{ \mathcal{G}(V, E) \mid |E| = m \right\}, \quad m = \sum_{i \in V} \hat{\mathbf{k}}_i^{\text{in}} = \sum_{i \in V} \hat{\mathbf{k}}_i^{\text{out}}, \quad (2)$$

with some probability measure P , such that the expected degree sequences of a realisation of X are fixed:

$$\mathbb{E}_P[\mathbf{k}^{\text{in}}(X)] = \hat{\mathbf{k}}^{\text{in}}, \quad \mathbb{E}_P[\mathbf{k}^{\text{out}}(X)] = \hat{\mathbf{k}}^{\text{out}}. \quad (3)$$

This set-up allows to map the model to an *urn problem* and thus to arrive to a closed-form probability distribution function for it.

Lemma 1 (Number of stub combinations). *The combinatorial matrix $\Xi \in \mathbb{N}^n \times \mathbb{N}^n$ given in Definition 4 encodes the numbers of out-stub and in-stub combinations for each pair of vertices, given degree sequences \mathbf{k}^{out} and \mathbf{k}^{in} .*

Proof. Let $\mathbf{k}_i^{\text{out}}$ be the out-degree of vertex i and \mathbf{k}_j^{in} the in-degree of vertex j . The number of out-stubs of a vertex corresponds to its out-degree. Similarly, the number of in-stubs of a vertex corresponds to its in-degree. Each one of the $\mathbf{k}_i^{\text{out}}$ out-stubs can be connected to all \mathbf{k}_j^{in} in-stubs. Hence, the total number of stub combinations between vertices i and j Ξ_{ij} is equal to $\mathbf{k}_i^{\text{out}}\mathbf{k}_j^{\text{in}}$. \square

We further introduce the concept of *induced random model*.

Definition 6 (Graph-induced random model). We say that the graph $\mathcal{G}(V, E)$ induces the random model X , if the quantities $(\hat{V}, \hat{\mathbf{k}}^{\text{in}}, \hat{\mathbf{k}}^{\text{out}})$ generating X are computed from \mathcal{G} . I.e., $\hat{V} = V$, $\hat{\mathbf{k}}^{\text{in}} = \mathbf{k}^{\text{in}}(\mathcal{G})$, and $\hat{\mathbf{k}}^{\text{out}} = \mathbf{k}^{\text{out}}(\mathcal{G})$.

Under this assumption, we can formulate the following theorem for the distribution of the soft configuration model X . To keep the notation simple, we will not distinguish between the $n \times n$ adjacency matrix \mathbf{A} and the vector of length n^2 obtained by stacking it by row or column. Similarly, we do the same for all other related $n \times n$ matrices.

Theorem 2. *Let $\mathcal{G}(V, E)$ be a directed graph with $n = |V|$ vertices and $m = |E|$ multi-edges. Let $\mathbf{k}^{\text{in}}(\mathcal{G}) \in \mathbb{N}^n$ and $\mathbf{k}^{\text{out}}(\mathcal{G}) \in \mathbb{N}^n$ be the vectors representing its in-degree and out-degree sequences. Let X be the soft configuration model induced by \mathcal{G} defined as in Definition 5. If the probability measure P depends only on the degree sequences in \mathcal{G} and the total number of multi-edges $m = |E|$ and all multi-edges are equiprobable, then X follows the multivariate hypergeometric distribution as in Eq. (4).*

Let $\mathbf{A} \in \mathbb{N}_0^n \times \mathbb{N}_0^n$ be an adjacency matrix and $\Xi \in \mathbb{N}^n \times \mathbb{N}^n$ be the combinatorial matrix induced by \mathcal{G} . Then the soft configuration model X is distributed as follows:

$$\Pr(X = \mathcal{G}) = \frac{\prod_{i,j \in V} \binom{\Xi_{ij}}{A_{ij}}}{\binom{M}{m}}, \quad (4)$$

where $M = \sum_{i,j \in V} \Xi_{ij}$ is the total number of stub combinations between all vertices.

Proof. We want to sample m multi-edges connecting any of the in- and out-stub pairs such that all such such multi-edges are equiprobable. According to Lemma 1, the total number of stubs combinations between any two vertices i, j is given by Ξ_{ij} . We can hence define the random graph model as follows. We sample m multi-edges without replacement from the multi-set of size $\sum_{i,j \in V} \Xi_{ij}$ that combines all the possible stub pairs combinations Ξ_{ij} between all pairs $i, j \in V$. We sample without replacement because we need to mimic the process of wiring stubs. Once a stub pair has been used it cannot be sampled again. We can view this model as an *urn problem* where the edges to be sampled are represented by balls in an urn. By representing the multi-edges connecting each different pair of vertices (i, j) as balls of a unique colour, we obtain an urn with a total of $M = \sum_{i,j \in V} \Xi_{ij}$ balls of $n^2 = |V \times V|$ different colours. With this, the sampling of a graph according to our model corresponds to drawing exactly m balls from this urn. Each adjacency matrix \mathbf{A} with $\sum_{i,j \in V} A_{ij} = m$ corresponds to one particular realisation drawn from this model. The probability to draw exactly $\mathbf{A} = \{A_{ij}\}_{i,j \in V}$ edges between each pair of vertices is given by the multivariate hypergeometric distribution. \square

From Theorem 2 we derive the following results, whose proofs follow directly from properties of the hypergeometric distribution.

Corollary 2.1. For each pair of vertices $i, j \in V$, the probability that X has exactly A_{ij} edges between i and j is given by the marginal distributions of the multivariate hypergeometric distribution, i.e.

$$\Pr(X_{ij} = A_{ij}) = \frac{\binom{\Xi_{ij}}{A_{ij}} \binom{M - \Xi_{ij}}{m - A_{ij}}}{\binom{M}{m}}. \quad (5)$$

Corollary 2.2. The expected in- and out-degree sequences of realisations of the directed soft configuration model X correspond to the respective degree sequences of the graph \mathcal{G} inducing X .

Proof. For each pair of vertices i, j we can calculate the expected number of multi-edges $\mathbb{E}[X_{ij}]$ as

$$\mathbb{E}[X_{ij}] = m \frac{\Xi_{ij}}{M} \quad (6)$$

Moreover, summing the rows and columns of matrix $\mathbb{E}[X_{ij}]$ and assuming directed graphs with self-loops we can calculate the expected in- or out-degrees of all vertices as

$$\begin{aligned} \mathbb{E}[k_j^{\text{in}}(X)] &= \sum_{i \in V} \mathbb{E}[X_{ij}] = m \frac{\sum_{i \in V} \hat{k}_i^{\text{out}} \hat{k}_j^{\text{in}}}{M} = \hat{k}_j^{\text{in}}, \\ \mathbb{E}[k_i^{\text{out}}(X)] &= \sum_{j \in V} \mathbb{E}[X_{ij}] = m \frac{\sum_{j \in V} \hat{k}_i^{\text{out}} \hat{k}_j^{\text{in}}}{M} = \hat{k}_i^{\text{out}}. \end{aligned} \quad (7)$$

Equation (7) confirms that the *expected* in- and out-degree sequence of realisations drawn from X corresponds to the degree sequence of the given graph \mathcal{G} . \square

2.1.2 Soft configuration model for undirected graphs

So far we have discussed the configuration model for directed graphs. Specifying the undirected case, on the other hand, requires some more efforts. The reason for this is that, under the assumptions described in the previous section, the undirected version of the soft configuration model is the degenerate case of its directed counterpart, where the direction of the multi-edges is ignored. In particular, this implies that the random vector corresponding to the undirected model has half the dimensions of the directed one, because any undirected multi-edge between two vertices i and j can either be generated as a directed multi-edge (i, j) or as a directed multi-edge (j, i) .

Definition 7 (Undirected soft configuration model). Let $\hat{\mathbf{k}} \in \mathbb{N}^n$ be a degree sequence and \hat{V} a set of n vertices. The soft configuration model X generated by $(\hat{V}, \hat{\mathbf{k}})$ is the $(n^2 + n)/2$ -dimensional random vector X defined on the probability space (S, P) , for the sample space

$$S = \left\{ \mathcal{G}(V, E) \mid |E| = m \right\}, \quad 2m = \sum_{i \in V} \hat{\mathbf{k}}_i, \quad (8)$$

with some probability measure P , such that the expected degree sequence of a realisation of X is fixed:

$$\mathbb{E}_P[\mathbf{k}(X)] = \hat{\mathbf{k}}. \quad (9)$$

At first sight, it would appear that the undirected soft configuration model can simply be obtained by restricting the directed model to $n(n + 1)/2$ components corresponding to the upper-triangle and the diagonal of the adjacency matrix. That is, we would sample m multi-edges among pairs $i \leq j \in V$ from the multi-set of stub combinations Ξ_{ij} as defined in Eq. (1). However, the resulting model does not satisfy Definition 7, because its expected degree sequence does not equal the degree sequence that induced it. To show this, we follow the same reasoning adopted in the proof of Corollary 2.2. The expected degree $\mathbb{E}[k_i(X)]$ of a vertex i is equivalent to

$$\mathbb{E}[k_i(X)] = \sum_{j \in V} \mathbb{E}[X_{ij}] = m \frac{\sum_{i \in V} \hat{\mathbf{k}}_i \hat{\mathbf{k}}_j}{\sum_{i \leq j \in V} \hat{\mathbf{k}}_i \hat{\mathbf{k}}_j} \neq \hat{\mathbf{k}}_i.$$

This approach is wrong because the total number of undirected stub combinations is larger than in the directed case. The reason for this is the symmetry in the process of wiring two stubs. Let $\hat{\mathbf{k}}_i$ be the degree of vertex i and $\hat{\mathbf{k}}_j$ the degree of vertex j . To form a multi-edge (i, j) , each one of the $\hat{\mathbf{k}}_i$ stubs of i can be connected to all $\hat{\mathbf{k}}_j$ stubs of j , and vice versa, each of the $\hat{\mathbf{k}}_j$ stubs of j can be connected to all $\hat{\mathbf{k}}_i$ stubs of i . Hence, the total number of combinations of stubs between vertices i and j equals to $\hat{\mathbf{k}}_i \hat{\mathbf{k}}_j + \hat{\mathbf{k}}_j \hat{\mathbf{k}}_i = 2\hat{\mathbf{k}}_i \hat{\mathbf{k}}_j$. As an equivalent to Lemma 1, we formalise this in the following lemma.

Lemma 3 (Undirected stubs combination count). *The total number of stub combinations between two vertices i, j in an undirected graph is given by:*

$$\begin{cases} 2\Xi_{ij} & \text{if } i \neq j, \\ \Xi_{ii} & \text{if } i = j, \end{cases} \quad (10)$$

We can now formulate the equivalent of Theorem 2 for the undirected case. The distribution underlying the undirected soft configuration model can then be computed analogously to Theorem 2 with the help of Lemma 3.

Theorem 4. Let $\mathcal{G}(V, E)$ be an undirected graph with $n = |V|$ vertices and $m = |E|$ multi-edges. Let $\mathbf{k} \in \mathbb{N}^n$ be the vector representing its degree sequence. Let X be the undirected soft configuration model induced by \mathcal{G} defined as in Definition 7. If the probability distribution underlying X depends only on the degree sequence of \mathcal{G} and the total number of multi-edges m , and all multi-edges are equiprobable, then X follows the multivariate hypergeometric distribution given in Eq. (11).

Let $\mathbf{A} \in \mathbb{N}_0^n \times \mathbb{N}_0^n$ be the symmetric adjacency matrix corresponding to an undirected graph, and $\Xi_{ij} = \mathbf{k}_i \mathbf{k}_j$ be the combinatorial matrix induced by \mathcal{G} . Then the undirected soft configuration model X is distributed as follows:

$$\Pr(X = \mathcal{G}) = \frac{\prod_{i < j \in V} \binom{2\Xi_{ij}}{A_{ij}} \prod_{l \in V} \binom{\Xi_{ll}}{A_{ll}/2}}{\binom{M}{m}}, \quad (11)$$

where $M = \sum_{i < j \in V} 2\Xi_{ij} + \sum_{l \in V} \Xi_{ll} = \sum_{i, j \in V} \Xi_{ij}$ is the total number of undirected stub combinations between all pair of vertices.

Proof. The proof follows the same reasoning of the proof of Theorem 2, accounting for the fact that the total number of stubs combinations is now given by Eq. (10). \square

Corollary 4.1. For each pair of vertices $i, j \in V$, the probability that X has exactly A_{ij} edges between i and j is given by the marginal distributions of the multivariate hypergeometric distribution in Eq. (11), i.e.,

$$\Pr(X_{ij} = A_{ij}) = \begin{cases} \binom{2\Xi_{ij}}{A_{ij}} \binom{M - 2\Xi_{ij}}{m - A_{ij}} \binom{M}{m}^{-1} & \text{for } i \neq j, \\ \binom{\Xi_{ij}}{A_{ij}/2} \binom{M - \Xi_{ij}}{m - A_{ij}/2} \binom{M}{m}^{-1} & \text{for } i = j. \end{cases} \quad (12)$$

Corollary 4.2. The expected degree sequence of realisations of X correspond to the respective degree sequences of the graph \mathcal{G} inducing X .

Proof. For each pair of vertices $i, j \in V$, the expected number of multi-edges $\mathbb{E}[X_{ij}]$ according to the hypergeometric distribution in Eq. (11) is expressed as

$$\mathbb{E}[X_{ij}] = 2m \frac{\Xi_{ij}}{M} \quad (13)$$

With this, we can write the expected degrees as

$$\mathbb{E}[k_j(X)] = \sum_{i \in V} \mathbb{E}[X_{ij}] = 2m \frac{\sum_{i \in V} \Xi_{ij}}{M} = 2m \frac{\sum_{i \in V} \hat{k}_i \hat{k}_j}{\sum_{i, j \in V} \hat{k}_i \hat{k}_j} = \hat{k}_j. \quad (14)$$

\square

2.1.3 Correspondence between directed and undirected models

In the previous two sections, we have formulated the soft configuration model for directed and undirected graphs independently of each other. As the reader recalls, we have motivated these models by the need of an analytically tractable analogy for the rewiring algorithm of the Molloy-Reed model [14]. This algorithm is the same in the directed and undirected case: select the first stub (outgoing, in the directed case), then select the second stub (incoming, in the directed case), create an edge by wiring these two stubs, and repeat the process until all the stubs are wired. Hence, we also show the correspondence between our directed and undirected formulations in this section.

We prove that the probability distribution of undirected graphs in the undirected soft configuration model given by Eq. (11) is a degenerate case of the directed model given by Eq. (4).

With the following definition we provide a projection from \mathbb{N}^{n^2} to $\mathbb{N}^{n(n+1)/2}$ that serves the purpose mapping a directed graph to its undirected equivalent, i.e., stripping the direction from its edges.

Definition 8 (Undirected projection). Let $\mathcal{G}^\dagger(V, E^\dagger)$ be a directed graph with adjacency matrix \mathbf{A}^\dagger . We define as *undirected projection* the map $\pi : \mathbb{N}^{n^2} \rightarrow \mathbb{N}^{n(n+1)/2}$ that maps \mathcal{G}^\dagger to the undirected graph $\mathcal{G}(V, E)$ with adjacency matrix $\mathbf{A} = \mathbf{A}^\dagger + \mathbf{A}^{\dagger T}$. We indicate with $\mathcal{G}^\dagger \hookrightarrow \mathcal{G}$ the fact that $\mathcal{G} = \pi(\mathcal{G}^\dagger)$.

According to Definition 8 there are different directed graphs that can be projected to the same undirected graph. At the same time, every undirected graph has at least one corresponding directed graph that can be projected to it, and for every directed graph there is at least one undirected graph to which it can be projected. These make the projection in Definition 8 surjective and not injective.

Similarly, we can define an undirected random graph model as the projection of a directed random graph model.

Definition 9. Let X^\dagger be a directed random graph model. With an abuse of notation, we use $X^{\dagger T}$ to refer to the transposition of the matrix representation of X^\dagger . We say that $X := X^\dagger + X^{\dagger T}$ is the undirected projection of X^\dagger if $\forall \mathcal{G}^\dagger$ in the sample space of X^\dagger exists a \mathcal{G} in the sample space of X such that the undirected projection $\pi(\mathcal{G}^\dagger)$ of \mathcal{G}^\dagger is \mathcal{G} . Furthermore, for every undirected graph \mathcal{G} in the sample space of X , $\exists \mathcal{G}^\dagger$ such that $\pi(\mathcal{G}^\dagger) = \mathcal{G}$. We indicate with $X^\dagger \hookrightarrow X$ the fact that X is the undirected projection of X^\dagger .

Note that according to Definition 8, the number of multi-edges m of \mathcal{G} equals the number of multi-edges of any directed \mathcal{G}^\dagger that projects to \mathcal{G} .

Finally we need to relate the distribution underlying a directed random graph model to the the distribution of its undirected projection. The following lemma serves this purpose.

Lemma 5. *Let \mathcal{G} be an undirected graph and X an undirected random graph model. Let X^\dagger be a directed random graph model such that $X^\dagger \hookrightarrow X$. The probability distribution of X , $\Pr(X = \mathcal{G})$, is given as:*

$$\Pr(X = \mathcal{G}) = \sum_{\mathcal{G}^\dagger \in \pi^{-1}(\mathcal{G})} \Pr(X^\dagger = \mathcal{G}^\dagger), \quad (15)$$

where the set $\pi^{-1}(\mathcal{G}) = \{\mathcal{G}^\dagger \mid \mathcal{G}^\dagger \hookrightarrow \mathcal{G}\}$ is the set of all directed graphs \mathcal{G}^\dagger that map to \mathcal{G} .

Proof. Let X^\dagger be a n^2 -dimensional random vector formalising a directed random graph model, such that $X^\dagger \hookrightarrow X$. For simplicity we index the elements of both random vectors as in the equivalent adjacency matrix notation. Let $X_{ij} = X_{ji}$ the ij -th element of X and $X_{ij}^\dagger, X_{ji}^\dagger$ the corresponding elements of X^\dagger .

According to Definition 9, X is the $n(n+1)/2$ -dimensional random vector defined as $X^\dagger + X^{\dagger T}$, where its each element ij is defined as $X_{ij} = X_{ij}^\dagger + X_{ji}^{\dagger T}$. The probability distribution of X , $\Pr(X = \mathcal{G}) = f_X(\mathcal{G})$ can be specified in terms of the probability distribution $f_{X^\dagger}(\mathcal{G}^\dagger) = \Pr(X^\dagger = \mathcal{G}^\dagger)$:

$$f_X(z) = f_X(\{z_{ij}\}_{ij}) = \sum \cdots \sum_{a_{ij}=0}^{z_{ij}} f_{X^\dagger}(\{z_{ij} - a_{ij}, a_{ij}\}_{ij,ji}) \quad (16)$$

The summation in Eq. (16) corresponds to the sum over the probabilities of all possible combinations of tuples $X_{ij}^\dagger, X_{ji}^\dagger$ which sum to A_{ij} for all indices ij . Hence, following Definition 8, this is equivalent to sum over all possible $\mathcal{G}^\dagger \hookrightarrow \mathcal{G}$. This proves the equivalence between Eq. (16) and Eq. (15) and thus, the lemma. \square

We can proceed showing that the undirected version of the soft-configuration model given in Theorem 2 is indeed equivalent to the model defined in Theorem 4. Theorem 6 stems from the fact that sampling an undirected edge between two vertices is equivalent to sampling a directed edge between the same pair of vertices in any of the two directions, and then stripping its direction information. The distribution underlying the undirected soft configuration model can then be computed with the help of Lemma 5.

Theorem 6. *Let \mathcal{G} be an undirected graph and X the undirected soft-configuration model. Let X^\dagger be the directed soft configuration model with combinatorial matrix with elements $\Xi_{ij} = \mathbf{k}_i \mathbf{k}_j$. The probability distribution of X is then given by Eq. (11).*

Proof. The distribution of X^\dagger is given by the hypergeometric distribution in Eq. (4). The model X^\dagger satisfies the conditions in Lemma 5 for the undirected soft configuration model X , because X^\dagger maps to the undirected soft configuration model X in accordance with Definition 9. Hence, we write the probability distribution underlying X as the sum of the probabilities of all corresponding directed graphs \mathcal{G}^\dagger under the directed soft configuration model X^\dagger .

$$\Pr(X = \mathcal{G}) = \sum_{\mathcal{G}^\dagger \in \pi^{-1}(\mathcal{G})} \Pr(X^\dagger = \mathcal{G}^\dagger) \quad (17)$$

$$= \sum_{\mathcal{G}^\dagger \in \pi^{-1}(\mathcal{G})} \binom{M}{m}^{-1} \prod_{i,j \in V} \binom{\Xi_{ij}}{A_{ij}^\dagger} \quad (18)$$

$$= \sum_{\mathcal{G}^\dagger \in \pi^{-1}(\mathcal{G})} \binom{M}{m}^{-1} \prod_{l \in V} \binom{\Xi_{ll}}{A_{ll}^\dagger} \prod_{i < j \in V} \binom{\Xi_{ij}}{A_{ij}^\dagger} \binom{\Xi_{ij}}{A_{ij} - A_{ij}^\dagger} \quad (19)$$

$$= \sum \dots \sum_{A_{ij}^\dagger=0}^{A_{ij}} \binom{M}{m}^{-1} \prod_{l \in V} \binom{\Xi_{ll}}{A_{ll}^\dagger} \prod_{i < j \in V} \binom{\Xi_{ij}}{A_{ij}^\dagger} \binom{\Xi_{ij}}{A_{ij} - A_{ij}^\dagger} \quad (20)$$

In Eq. (20) we have $n(n-1)/2$ summations for all A_{ij}^\dagger , $i < j$, which are the decomposition of the summation in Eq. (19). Then, we can swap the summations and multiplications in Eq. (20), which leads to

$$\Pr(\mathbf{A}) = \binom{M}{m}^{-1} \prod_{l \in V} \binom{\Xi_{ll}}{A_{ll}^\dagger} \prod_{i < j \in V} \sum_{A_{ij}^\dagger=0}^{A_{ij}} \binom{\Xi_{ij}}{A_{ij}^\dagger} \binom{\Xi_{ij}}{A_{ij} - A_{ij}^\dagger}. \quad (21)$$

From Vandermonde's identity, which states

$$\sum_{a=0}^{A_{ij}} \binom{\Xi_{ij}}{a} \binom{2\Xi_{ij} - \Xi_{ij}}{A_{ij} - a} = \binom{2\Xi_{ij}}{A_{ij}}, \quad (22)$$

and from the fact that $A_{ii}^\dagger = A_{ii}/2$, $\forall i \in V$, it follows that Eq. (21) is equivalent to Eq. (11). \square

In the two sections above we have provided a parsimonious formulation of the soft-configuration model in terms of an *hypergeometric ensemble*. The ensemble provides a random graph model formulation for directed and undirected graphs alike, in which (i) the expected in- and out-degree sequences are fixed, and (ii) multi-edges between these vertices with fixed expected degrees are formed uniformly at random. More precisely, the probability for a particular pair of vertices to be connected by an edge is only influenced by combinatorial effects, and thus only depends on the degrees of the vertices and the total number of edges.

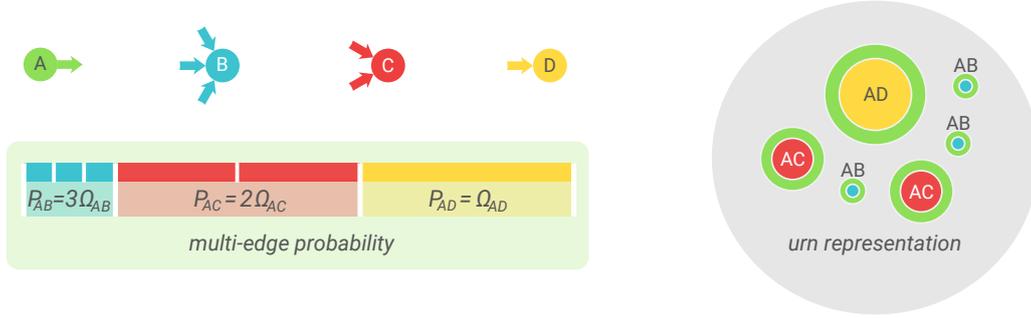


Figure 2: The effect of *edge propensities* on the configuration model. Differently from the standard configuration model, here the stubs are not sampled uniformly at random (cf. Fig. 1). Given an out-stub, each in-stub is characterised by a propensity Ω_{ij} of being chosen. As a result, the probability of wiring the out-stub (A, \cdot) to the vertex D is larger than that of B due to a very large edge propensity Ω_{AD} , even though vertex B has three times more in-stubs than vertex D .

2.2 Generalised Hypergeometric Ensemble of Graphs

As described above, the hypergeometric ensemble of random graphs samples edges uniformly at random from the urn containing all possible edges. However, such uniform sampling of edges is generally not enough to describe real graphs. In fact, in empirical graphs edge probabilities do not only depend on the activity of vertices (represented by their degrees) but also on other characteristics. Examples of such characteristics are vertex labels that lead to observable group structures, distances between vertices, etc. Below we introduce how such influences on the sampling probabilities of edges can be encoded in a random graph model by means of a dyadic property we call *edge propensity*.

First, we introduce the concept of edge propensity. Given two dyads (i, j) and $(k, l) \in V \times V$ where $\Xi_{ij} = \Xi_{kl}$, according to the soft-configuration model of Section 2.1 the probabilities of sampling one multi-edge between (i, j) and (k, l) are equal, leading to odds-ratio of 1 between the two pairs of vertices. Instead, we generalise the model in a way that fixing arbitrary odds-ratios is possible. That is, we define the aforementioned edge propensities such that the ratio between them is the odds-ratio of sampling one multi-edge between the two corresponding vertex pairs, all else being equal. We use edge propensities to bias the sampling probability of each multi-edge, as illustrated in Fig. 2. This way, the probability of the number of multi-edges between a pair of vertices depends on both the degrees of the vertices and their edge propensity.

We encode edge propensities in a matrix $\mathbf{\Omega}$ defined as follows.

Definition 10 (Propensity matrix). Let $\mathbf{\Omega} = (\Omega_{ij})_{i,j \in V} \in \mathbb{R}^{|V| \times |V|}$ be a $n \times n$ matrix where n is the number of vertices. Let (i, j) and $(k, l) \in V \times V$. Let $\omega_{ij,kl}$ the odds-ratio of sampling one multi-edge between (i, j) instead of (k, l) . The entries Ω_{ij} and Ω_{kl} of the propensity matrix $\mathbf{\Omega}$ are then defined such that $\Omega_{ij}/\Omega_{kl} = \omega_{ij,kl}$. This implies that the propensity matrix $\mathbf{\Omega}$ is defined up to a constant, as multiplying $\mathbf{\Omega}$ with any constant preserves the specified odds-ratios.

Now, we define a random graph model that combines the degree-related combinatorial effects, i.e., the configuration model, and the newly introduced edge propensities. We do so by using the propensities to *bias* the sampling process described in Section 2.1. In the urn model analogy, such biased sampling implies that the probability of drawing a certain number of balls of a given colour (i.e., multi-edges between the corresponding pair of vertices) depends both on their number and their size, as illustrated in Fig. 2. The probability distribution resulting from such a biased sampling process is given by the multivariate Wallenius' non-central hypergeometric distribution [18, 11].

Theorem 7. Let $\mathcal{G}(V, E)$ be a directed graph with $n = |V|$ vertices and $m = |E|$ edges. Under the assumptions introduced above, the generalised hypergeometric ensemble of graphs (GHypEG) X induced by \mathcal{G} and a given propensity matrix $\mathbf{\Omega}$ follows the multivariate Wallenius' non-central hypergeometric distribution given in Eq. (23).

Let $\mathbf{A} \in \mathbb{N}^n \times \mathbb{N}^n$ be the adjacency matrix associated with \mathcal{G} and $\mathbf{\Xi} \in \mathbb{N}^n \times \mathbb{N}^n$ be its combinatorial matrix defined in Eq. (1). Then the GHypEG defined by $\mathbf{\Xi}$ and $\mathbf{\Omega}$, \mathcal{G} is distributed as follows:

$$\Pr(X = \mathbf{A}) = \left[\prod_{i,j \in V} \binom{\Xi_{ij}}{A_{ij}} \right] \int_0^1 \prod_{i,j \in V} \left(1 - z \frac{\Omega_{ij}}{S_{\mathbf{\Omega}}} \right)^{A_{ij}} dz \quad (23)$$

with

$$S_{\mathbf{\Omega}} = \sum_{i,j \in V} \Omega_{ij} (\Xi_{ij} - A_{ij}). \quad (24)$$

The distribution describing the biased sampling from an urn is a generalisation of the multivariate hypergeometric distribution. The proof of Theorem 7 follows from the fact that, when the sampling is performed without replacement with given relative odds, this sampling process corresponds to the multivariate Wallenius' non-central hypergeometric distribution. Details of this derivation can be found in [18] for the univariate case, and in [7, 6] for the multivariate case. A thorough review of non-central hypergeometric distributions has been done in [12].

The next two corollaries directly follow from properties of Wallenius' non-central hypergeometric distribution.

Corollary 7.1. *For each pair of vertices $i, j \in V$, the probability to draw exactly A_{ij} edges between i and j is given by the marginal distributions of the multivariate Wallenius' non-central hypergeometric distribution, i.e.,*

$$\Pr(X_{ij} = A_{ij}) = \binom{\Xi_{ij}}{A_{ij}} \binom{M - \Xi_{ij}}{m - A_{ij}} \cdot \int_0^1 \left[\left(1 - z \frac{\Omega_{ij}}{s\bar{\Omega}}\right)^{A_{ij}} \left(1 - z \frac{\bar{\Omega}_{ij}}{s\bar{\Omega}}\right)^{m - A_{ij}} \right] dz \quad (25)$$

where

$$\bar{\Omega}_{ij} = \frac{\sum_{(l,m) \in (V \times V) \setminus (i,j)} \Xi_{lm} \Omega_{lm}}{(M - \Xi_{ij})}. \quad (26)$$

Corollary 7.2. *The entries of the expected adjacency matrix $\mathbb{E}[X_{ij}]$ can be obtained by solving the following system of equations:*

$$\left(1 - \frac{\mathbb{E}[X_{11}]}{\Xi_{11}}\right)^{\frac{1}{\bar{\Omega}_{11}}} = \left(1 - \frac{\mathbb{E}[X_{12}]}{\Xi_{12}}\right)^{\frac{1}{\bar{\Omega}_{12}}} = \dots \quad (27)$$

with the constraint $\sum_{i,j \in V} \mathbb{E}[X_{ij}] = m$.

The undirected formulation of the GHypEG follows the same reasoning of Theorem 4, with the addition of a symmetric propensity matrix. That is, the distribution of the undirected generalised hypergeometric ensemble is hence given by a Wallenius' distribution similar to Eq. (23), but corresponding to the upper triangular part of the matrices (i.e., for $i \leq j$) with Ξ defined according to Lemma 3.

2.3 Estimation of the propensity matrix

In this final section, we show how to define the propensity matrix Ω such that the expected graph $\mathbb{E}[X]$ from the model X defined by Ω coincides with an arbitrary graph G . By doing so, we create a random model *centered* around the inducing graph. The result is described in the following corollary, which follows from the properties of Wallenius' non-central hypergeometric distribution.

Corollary 7.3. *Let $\mathcal{G}(V, E)$ be a graph with $n = |V|$ vertices and \mathbf{A} its adjacency matrix. Let Ω be a $n \times n$ propensity matrix, characterised by elements Ω_{ij} with $i, j \in V$. Then, \mathcal{G} coincides with the expectation $\mathbb{E}[X]$ of the GHypEG X induced by \mathcal{G} and Ω if and only if the following relation holds.*

$$\forall c \in \mathbb{R}^- \quad \Omega_{ij} = \frac{1}{c} \log\left(1 - A_{ij}/\Xi_{ij}\right) \quad \forall i, j \in V, \quad (28)$$

Where Ξ is the combinatorial matrix associated with X and Ξ_{ij} its elements.

Proof. Equation (28) follows directly from Corollary 7.1. In particular, when solving Eq. (27) for Ω with the assumption $\mathbb{E}[X] = \mathbf{A}$ we obtain the following system of $|V|^2$ equations (in the case of a directed graph with self-loops) for $|V|^2 + 1$ variables.

$$\begin{cases} \left(1 - \frac{\mathbb{E}[X_{11}]}{\Xi_{11}}\right)^{\frac{1}{\Omega_{11}}} = C \\ \left(1 - \frac{\mathbb{E}[X_{12}]}{\Xi_{12}}\right)^{\frac{1}{\Omega_{12}}} = C \\ \vdots \end{cases} \quad (29)$$

The solution of this system is Eq. (28). □

A wide range of statistical patterns that go beyond degree effects can be encoded in the graph model by specifying the matrix Ω of edge propensities. The encoding and fitting techniques of such arbitrary propensity matrices are beyond the scope of this article, and will not be discussed here. We refer to [4] for a general method to fit external dyadic data as propensities.

Finally, we show that the soft configuration model of Section 2.1 is a special case of generalised hypergeometric graph models. The soft configuration model described in Theorem 2 can be in fact recovered from the generalised model by setting all entries in the propensity matrix to the same value. By doing so, the odds-ratio between the propensities for any pair of vertices is 1, and the edge sampling process is not biased. Thus, the probability distribution of the model reduces to a function of the degree sequences and the number of multi-edges sampled.

Theorem 8. *Let $\Omega \equiv \text{const}$. The corresponding GHypEG coincides with the soft configuration model in Eq. (2) induced by the same graph.*

Proof. For the special case of a uniform edge propensity matrix $\Omega \equiv \text{const}$, which corresponds to an unbiased sampling of edges, for the integral in Eq. (25) we have

$$\int_0^1 \left(1 - z^{\frac{1}{M-m}}\right)^m dz = \binom{M}{m}^{-1}. \quad (30)$$

Plugging this result in Eq. (23) we thus recover Eq. (4) for the unbiased case, i.e. where all edge propensities are identical. □

3 Final remarks

We have proposed a novel approach for studying random graphs in terms of an urn problem. By doing so, we have arrived at an analytically tractable formulation of the

widely used configuration model of random graphs. Furthermore, we have expanded the configuration model to the *generalised hypergeometric ensemble*, which is a whole new class of models that can incorporate arbitrary dyadic biases in the probabilities of sampling edges. Importantly, the analytical tractability of the generalised hypergeometric ensembles allows for robust model selection and hypothesis testing of various topological patterns in empirical data [5]. Moreover, one can perform a *multiplex network regression* [4] based on our model, in order to find the weighted combination of multiple relational layers that best describes the observed multi-graph. For instance, one can study how different social phenomena and environmental factors influence the topological patterns in repeated social interactions. The ability to analytically perform such statistical analysis opens new possibilities for the study of complex interaction data, which pervades many scientific disciplines.

Aknowledgements

The authors thank Frank Schweitzer for his support and valuable comments, and Ingo Scholtes, Pavlin Mavrodiev, and Christian Zingg for the useful discussions.

References

- [1] W. AIELLO, F. CHUNG, AND L. LU, *A random graph model for massive graphs*, in Proceedings of the 32nd annual ACM symposium on Theory of computing - STOC '00, ACM Press, 2000, pp. 171–180.
- [2] A.-L. BARABÁSI AND R. ALBERT, *Emergence of scaling in random networks*, science, 286 (1999), pp. 509–512.
- [3] E. A. BENDER AND E. CANFIELD, *The asymptotic number of labeled graphs with given degree sequences*, Journal of Combinatorial Theory, A, 24 (1978), pp. 296–307.
- [4] G. CASIRAGHI, *Multiplex Network Regression: How do relations drive interactions?*, arXiv preprint arXiv:1702.02048, (2017).
- [5] G. CASIRAGHI, V. NANUMYAN, I. SCHOLTES, AND F. SCHWEITZER, *Generalized Hypergeometric Ensembles: Statistical Hypothesis Testing in Complex Networks*, arXiv preprint arXiv:1607.02441, (2016).

- [6] J. CHESSON, *A non-central multivariate hypergeometric distribution arising from biased sampling with application to selective predation*, *Journal of Applied Probability*, (1976), pp. 795–797.
- [7] —, *Measuring Preference in Selective Predation*, *Ecology*, 59 (1978), pp. 211–215.
- [8] F. CHUNG AND L. LU, *Connected Components in Random Graphs with Given Expected Degree Sequences*, *Annals of Combinatorics*, 6 (2002), pp. 125–145.
- [9] F. CHUNG AND L. LU, *The average distances in random graphs with given expected degrees*, *Proceedings of the National Academy of Sciences*, 99 (2002), pp. 15879–15882.
- [10] P. ERDÖS AND A. RÉNYI, *On random graphs I*, *Publ. Math. Debrecen*, 6 (1959), pp. 290–297.
- [11] A. FOG, *Calculation Methods for Wallenius’ Noncentral Hypergeometric Distribution*, *Communications in Statistics - Simulation and Computation*, 37 (2008), pp. 258–273.
- [12] —, *Sampling Methods for Wallenius’ and Fisher’s Noncentral Hypergeometric Distributions*, *Communications in Statistics - Simulation and Computation*, 37 (2008), pp. 241–257.
- [13] B. A. HUBERMAN AND L. A. ADAMIC, *Growth dynamics of the World-Wide Web*, *Nature*, 401 (1999), pp. 131–131.
- [14] M. MOLLOY AND B. REED, *A critical point for random graphs with a given degree sequence*, *Random Structures & Algorithms*, 6 (1995), pp. 161–180.
- [15] —, *The Size of the Giant Component of a Random Graph with a Given Degree Sequence*, *Combinatorics, Probability and Computing*, 7 (1998), pp. 295–305.
- [16] M. E. J. NEWMAN, *Modularity and community structure in networks*, *Proceedings of the National Academy of Sciences*, 103 (2006), pp. 8577–8582.
- [17] M. E. J. NEWMAN AND J. PARK, *Why social networks are different from other types of networks*, *Phys. Rev. E*, 68 (2003), p. 036122.
- [18] K. T. WALLENIOUS, *Biased Sampling: the Noncentral Hypergeometric Probability Distribution*, ph.d. thesis, Stanford University, 1963.