

Entropies of tailored random graph ensembles: bipartite graphs, generalized degrees, and node neighbourhoods

E S Roberts^{1,2} and A C C Coolen^{1,3}

¹Institute for Mathematical and Molecular Biomedicine, King's College London, Hodgkin Building, London SE1 1UL, UK

²Randall Division of Cell and Molecular Biophysics, King's College London, New Hunts House, London SE1 1UL, UK

³London Institute for Mathematical Sciences, 35a South St, Mayfair, London W1K 2XF, UK

E-mail: ekaterina.roberts@kcl.ac.uk and ton.coolen@kcl.ac.uk

Received 23 April 2014, revised 20 August 2014

Accepted for publication 20 August 2014

Published 9 October 2014

Abstract

We calculate explicit formulae for the Shannon entropies of several families of tailored random graph ensembles for which no such formulae were as yet available, in leading orders in the system size. These include bipartite graph ensembles with imposed (and possibly distinct) degree distributions for the two node sets, graph ensembles constrained by specified node neighbourhood distributions, and graph ensembles constrained by specified generalized degree distributions.

Keywords: random graphs, networks, entropy, generalized degrees, bipartite graphs

PACS numbers: 89.70.Cf, 89.75.Fb, 64.60.aq

(Some figures may appear in colour only in the online journal)

1. Introduction

Networks are powerful and popular tools for characterizing large and complex interacting particle systems. They have become extremely valuable in physics, biology, computer science, economics, and the social sciences. One approach is to quantify the implications of having topological patterns in networks and graphs, by viewing these patterns as constraints

on a random graph ensemble. This provides a way to measure and compare topological features from the rational point of view of whether they are present in a large or small number of possible networks. Precise definitions of random graph ensembles with controlled topological characteristics also allow us to generate systematically graphs and networks which are tailored to have features in common with those observed in a given application domain, either for the purpose of statistical mechanical process modelling or to serve as ‘null models’ against which to test the importance of observations in real-world networks.

A previous paper [1] considered tailored random graph ensembles with controlled degree distribution and degree–degree correlations; the more recent [2] covered the case of directed networks. In each case, the strategy is to calculate the Shannon entropy, from which we can deduce the effective number of graphs in the ensemble. Related quantities such as complexity of typical graphs from the ensemble and information-theoretic distances between graphs naturally follow from the entropy, or can be calculated using similar methods.

In this paper we calculate, in leading order, the Shannon entropies of three as yet unsolved families of random graph ensembles, constrained by three different conditions: a bipartite constraint with imposed degree distributions in the two nodes sets, a neighbourhood distribution (where the neighbourhood of a node is defined as its own degree, plus the degree values of the nodes connected to it), and an imposed generalized degree distribution. These are each interesting in their own right as stand-alone results, and turn out to be closely linked. The first two cases can be resolved exactly, and give practical analytical expressions. The generalized degree case was already partially studied in [3], with only limited success, and here we require a plausible but as yet unproven conjecture to find an explicit formula for the entropy.

The generalized degrees concept appears in the literature in various forms. For example, the authors of [4], measured the number of direct neighbours s of a subset of t nodes. They derive conditions based on their definition of general degrees which can ensure that (for some given m and d) there are at least m internally disjoint paths of length at most d . The diameter of the network is an obvious corollary—the smallest d corresponding to $m \geq 1$. These results can be applied to questions of robustness of networks. The authors of [5] studied the spectral density of random graphs with hierarchically constrained topologies. This includes consideration of generalized degrees, as well as more general community structures. Using the replica method, in a similar way to [3], they achieve a form analogous to equation (39). They proceed numerically from that point, hence our approach to an analytical solution presented in equation (48) is entirely novel.

2. Definitions and notation

We consider ensembles of directed and nondirected random graphs. Each graph is defined by its adjacency matrix $c = \{c_{ij}\}$, with $i, j \in \{1, \dots, N\}$ and with $c_{ij} \in \{0, 1\}$ for all (i, j) . Two nodes i and j are connected by a directed link $j \rightarrow i$ if and only if $c_{ij} = 1$. We put $c_{ii} = 0$ for all i . In nondirected graphs one has $c_{ij} = c_{ji}$ for all (i, j) , so c is symmetric. The degree of a node i in a nondirected graph is the number of its neighbours, $k_i = \sum_j c_{ij}$. In directed graphs we distinguish between in- and out-degrees, $k_i^{\text{in}} = \sum_j c_{ij}$ and $k_i^{\text{out}} = \sum_j c_{ji}$. They count the number of in- and out-bound links at a node i . A bipartite graph is one where the nodes can be divided into two disjoint sets, such that $c_{ij} = 0$ for all i and j that belong to the same set.

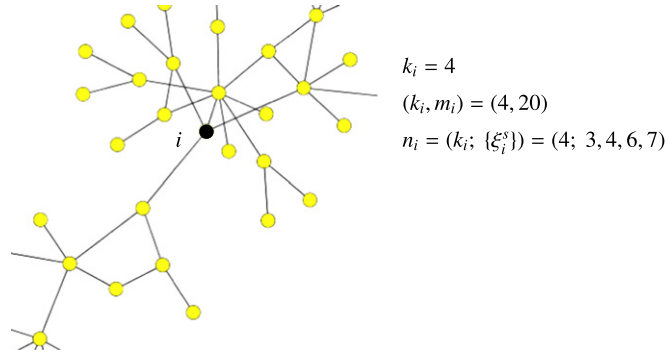


Figure 1. Illustration of our definitions of local topological characteristics in nondirected graphs. At the minimal level one specifies for each node i (black vertex in the picture) only the degree $k_i = |\partial_i| = \sum_j c_{ij}$ (the number of its neighbours). At the next level of detail one provides for each node the generalized degree (k_i, m_i) , in which $m_i = \sum_{j \in \partial_i} k_j = \sum_j c_{ij} k_j$ is the number of length-two paths starting in i . This is then generalized to include the actual degrees in the set ∂_i , by giving $n_i = (k_i; \{\xi_i^s\})$ (the ‘local neighbourhood’), in which the k_i integers $\{\xi_i^s\}$ give the degrees of the nodes connected to i . To avoid ambiguities we adopt the ranking convention $\xi_i^1 \leq \xi_i^2 \leq \dots \leq \xi_i^{k_i}$. Note that $m_i = \sum_{j \in \partial_i} k_j = \sum_{s=1}^{k_i} \xi_i^s$.

We define the set of neighbours of a node i in a nondirected graph as $\partial_i = \{j \mid c_{ij} = 1\}$. Hence $k_i = |\partial_i|$. To characterize a graph’s topology near i in more detail we can define the generalized degree of i as the pair (k_i, m_i) , where $m_i = \sum_j c_{ij} k_j$ counts the number of length-two paths starting in i . The concept of a generalized degree is discussed in [6]. Even more information is contained in the *local neighbourhood*

$$n_i = (k_i; \{\xi_i^s\}), \tag{1}$$

in which the ordered integers $\{\xi_i^s\}$ give the degrees of the k_i neighbours $j \in \partial_i$. See also figure 1. Since $m_i = \sum_{s \leq k_i} \xi_i^s$, the neighbourhood n_i provides more granular information that complements that in the generalized degree (k_i, m_i) . We will use bold symbols when local topological parameters are defined for every node in a network, e.g. $\mathbf{k} = (k_1, \dots, k_N)$ and $\mathbf{n} = ((k_1; \{\xi_1^s\}), \dots, (k_N; \{\xi_N^s\}))$. Generalization to directed graphs is straightforward. Here $\partial_i = \{j \mid c_{ij} + c_{ji} > 0\}$, and the local neighbourhood would be defined as $n_i = (\vec{k}_i; \{\vec{\xi}_i^s\})$ with the k_i pairs $\vec{\xi}_i^s = (k^{s,\text{in}}, k^{s,\text{out}})$ now giving both the in- and out-degrees of the neighbours of i .

Our tailored random graph ensembles will be of the following form, involving N built-in local (site specific) topological constraints of the type discussed above, which we will for now write generically as $X_i(\mathbf{c})$, and with the usual abbreviation $\delta_{a,b} = \prod_i \delta_{a_i,b_i}$

$$\begin{aligned} p(\mathbf{c}) &= \sum_{\mathbf{X}} p(\mathbf{X}) p(\mathbf{c}|\mathbf{X}) & p(\mathbf{X}) &= \prod_i p(X_i), \\ p(\mathbf{c}|\mathbf{X}) &= Z^{-1}(\mathbf{X}) \delta_{\mathbf{X},\mathbf{X}(\mathbf{c})}, & Z(\mathbf{X}) &= \sum_{\mathbf{c}} \delta_{\mathbf{X},\mathbf{X}(\mathbf{c})}. \end{aligned} \tag{2}$$

The values X_i for the local features are for each i drawn randomly and independently from $p(X)$, after which one generates a graph \mathbf{c} randomly and with uniform probabilities from the

set of graphs that satisfy the N demands $X_i(\mathbf{c}) = X_i$. The empirical distribution $p(\mathbf{X}|\mathbf{c}) = N^{-1} \sum_i \delta_{X_i, X_i(\mathbf{c})}$ of local features will be random, but the law of large numbers ensures that for $N \rightarrow \infty$ it will converge to the chosen $p(X)$ in (2) for any graph realization, and the above definitions guarantee that its ensemble average will be identical to $p(X)$ for any N

$$\sum_{\mathbf{c}} p(\mathbf{c}) p(\mathbf{X}|\mathbf{c}) = \frac{1}{N} \sum_i \sum_X p(X) \sum_{\mathbf{c}} \frac{\delta_{X_i, X_i(\mathbf{c})}}{Z(\mathbf{X})} \delta_{X_i, X_i(\mathbf{c})} = p(X). \quad (3)$$

If we aim to impose upon our graphs only a degree distribution we choose $X_i(\mathbf{c}) = k_i(\mathbf{c})$. Building in a distribution of generalized degrees corresponds to $X_i(\mathbf{c}) = (k_i(\mathbf{c}), m_i(\mathbf{c}))$. If we seek to prescribe the distribution of all local neighbourhoods (1) we choose $X_i(\mathbf{c}) = n_i(\mathbf{c})$.

A further quantity which will play a role in subsequent calculations is the joint degree distribution of connected nodes. For nondirected graphs it is defined as

$$W(k, k'|\mathbf{c}) = \frac{\sum_{ij} c_{ij} \delta_{k, k_i} \delta_{k', k_j}}{\sum_{ij} c_{ij}} \quad (4)$$

and its average over the ensemble (2) is given by

$$W(k, k') = \sum_X p(X) \sum_{\mathbf{c}} W(k, k'|\mathbf{c}) \frac{\delta_{X, X(\mathbf{c})}}{Z(\mathbf{X})}. \quad (5)$$

In this paper we study the leading orders in the system size N of the Shannon entropy per node of the above tailored random graph ensembles (2), from which the effective number of graphs with the prescribed distribution $p(X)$ of features follows as $\mathcal{N} = \exp(NS)$

$$\begin{aligned} S &= -\frac{1}{N} \sum_{\mathbf{c}} p(\mathbf{c}) \log p(\mathbf{c}) \\ &= -\frac{1}{N} \sum_X \frac{\prod_i p(X_i)}{Z(\mathbf{X})} \sum_{\mathbf{c}} \delta_{X, X(\mathbf{c})} \log \left[\frac{\prod_{X'} p(X'_j)}{\sum_{X'} \frac{Z(X')^j}{Z(\mathbf{X})} \delta_{X', X(\mathbf{c})}} \right] \\ &= -\frac{1}{N} \sum_X \frac{\prod_i p(X_i)}{Z(\mathbf{X})} \sum_{\mathbf{c}} \delta_{X, X(\mathbf{c})} \log \left[\frac{\prod_j p(X_j)}{Z(\mathbf{X})} \right] \\ &= \sum_X p(X) S(X) - \sum_X p(X) \log p(X) \end{aligned} \quad (6)$$

with

$$S(X) = \frac{1}{N} \log Z(\mathbf{X}) = \frac{1}{N} \log \sum_{\mathbf{c}} \delta_{X, X(\mathbf{c})}. \quad (7)$$

The core of the entropy calculation is determining the leading orders in N of $S(\mathbf{X})$, which is the Shannon entropy per node of the ensemble $p(\mathbf{c}|\mathbf{X})$ in which all node-specific values $\mathbf{X} = (X_1, \dots, X_N)$ are constrained. For $p(X) = p(k)$ this calculation has already been done in [1, 2]. For $p(X) = p(k, m)$ it has only partly been done [3]. Here we investigate the relation

between the entropies of the $p(k)$ and $p(k, m)$ ensembles and the entropy of the ensemble in which the distribution $p(n)$ of local neighbourhoods (1) is imposed.

3. Building blocks of the entropy calculations

3.1. Relations between feature distributions for nondirected graphs

Since the generalized degrees (k_i, m_i) can be calculated from the local neighbourhoods (1) for any graph \mathbf{c} , it is clear that the empirical distribution $p(k, m|\mathbf{c}) = N^{-1} \sum_i \delta_{k,k_i(\mathbf{c})} \delta_{m,m_i(\mathbf{c})}$ for any graph can be calculated from the empirical neighbourhood distribution $p(n|\mathbf{c}) = N^{-1} \sum_i \delta_{n,n_i(\mathbf{c})}$. If we denote with $k(n)$ the central degree k in $n = (k; \{\xi^s\})$, we indeed obtain

$$p(k, m|\mathbf{c}) = \frac{1}{N} \sum_i \delta_{k,k_i(\mathbf{c})} \delta_{m,m_i(\mathbf{c})} \sum_n \delta_{n,n_i} = \sum_n p(n) \delta_{k,k(n)} \delta_{m, \sum_{s \leq k(n)} \xi^s}. \quad (8)$$

Less trivial is the statement that also the distribution $W(k, k'|\mathbf{c})$ of (4) can be written in terms of $p(n|\mathbf{c})$. Using $\sum_{ij} c_{ij} = N\bar{k}(\mathbf{c})$, with $\bar{k}(\mathbf{c}) = N^{-1} \sum_i k_i(\mathbf{c})$ we obtain

$$\begin{aligned} W(k, k'|\mathbf{c}) &= \frac{\sum_i \delta_{k,k_i(\mathbf{c})} \sum_{j \in \partial_i} \delta_{k',k_j(\mathbf{c})}}{N \sum_n p(n|\mathbf{c}) k(n)} = \frac{\sum_i \sum_n \delta_{n,n_i(\mathbf{c})} \delta_{k,k(n)} \sum_{s \leq k(n)} \delta_{k',\xi^s}}{N \sum_n p(n|\mathbf{c}) k(n)} \\ &= \frac{\sum_n p(n|\mathbf{c}) \delta_{k,k(n)} \sum_{s \leq k(n)} \delta_{k',\xi^s}}{\sum_n p(n|\mathbf{c}) k(n)}. \end{aligned} \quad (9)$$

Given the symmetry of $W(k, k'|\mathbf{c})$ under permutation of k and k' we then also have

$$W(k, k'|\mathbf{c}) = \frac{\sum_n p(n|\mathbf{c}) \delta_{k',k(n)} \sum_{s \leq k(n)} \delta_{k,\xi^s}}{\sum_n p(n|\mathbf{c}) k(n)}. \quad (10)$$

The converse of the above statements is not true. One cannot calculate the neighbourhood distribution $p(n|\mathbf{c})$ from $p(k, m|\mathbf{c})$ or from $W(k, k'|\mathbf{c})$ (or both). Note that by definition (and since \mathbf{c} is nondirected) we always have $W(k, k'|\mathbf{c}) = W(k', k|\mathbf{c})$.

3.2. Decomposition of graphs into directed degree-regular subgraphs

Any nondirected graph \mathbf{c} can always be decomposed uniquely into a collection of non-overlapping N -node subgraphs $\beta^{kk'}$, with $k, k' \in \mathbb{N}$, which share the nodes $\{1, \dots, N\}$ of \mathbf{c} but not all of the links. These subgraphs are defined for each (k, k') by the adjacency matrices

$$\beta_{ij}^{kk'} = c_{ij} \delta_{k,k_i(\mathbf{c})} \delta_{k',k_j(\mathbf{c})}. \quad (11)$$

Each graph $\beta^{kk'}$ contains those links in \mathbf{c} that go from a node with degree k' to a node with degree k . Clearly, all graphs $\beta^{kk'}$ follow uniquely from \mathbf{c} via (11). The converse uniqueness of \mathbf{c} , given the matrices $\beta^{kk'}$, is a consequence of the simple identity

$$c_{ij} = c_{ij} \sum_{kk' \geq 0} \delta_{k,k_i(\mathbf{c})} \delta_{k',k_j(\mathbf{c})} = \sum_{kk' \geq 0} \delta_{k,k_i(\mathbf{c})} \delta_{k',k_j(\mathbf{c})} c_{ij} = \sum_{kk' \geq 0} \beta_{ij}^{kk'}. \quad (12)$$

The graph $\beta^{kk'}$ is directed if $k \neq k'$, and nondirected if $k = k'$. From the symmetry of \mathbf{c} it follows moreover that $\beta_{ji}^{kk'} = \beta_{ij}^{k'k}$ for all (i, j, k, k') , so $\beta^{kk'}$ is specified in full by $\beta^{kk'}$.

Although each $\beta^{kk'}$ is an N -node graph, most of the nodes in $\beta^{kk'}$ will be isolated: all nodes whose degrees in the original graph \mathbf{c} were neither k nor k' will have degree zero in $\beta^{kk'}$.

We now inspect the degree statistics of the decomposition graphs $\beta^{kk'}$, and their relation with the structural features of \mathbf{c} . If $k \neq k'$ we find for the remaining degrees in $\beta^{kk'}$:

$$k_i(\mathbf{c}) = k: \quad k_i^{\text{in}}(\beta^{kk'}) = \sum_{j \in \partial_i} \delta_{k',k_j(\mathbf{c})}, \quad k_i^{\text{out}}(\beta^{kk'}) = 0, \quad (13)$$

$$k_j(\mathbf{c}) = k': \quad k_j^{\text{out}}(\beta^{kk'}) = \sum_{i \in \partial_j} \delta_{k,k_i(\mathbf{c})}, \quad k_j^{\text{in}}(\beta^{kk'}) = 0. \quad (14)$$

Hence the joint in–out degree distribution of $\beta^{kk'}$ can be written in terms of the empirical distribution of neighbourhoods of \mathbf{c} , viz. $p(n|\mathbf{c}) = N^{-1} \sum_i \delta_{n,n_i(\mathbf{c})}$ with $n = (k; \{\xi^s\})$

$$\begin{aligned} p^{kk'}(q^{\text{in}}, q^{\text{out}}) &= \frac{1}{N} \sum_i \delta_{q^{\text{in}},k_i^{\text{in}}(\beta^{kk'})} \delta_{q^{\text{out}},k_i^{\text{out}}(\beta^{kk'})} \\ &= \frac{1}{N} \sum_i \delta_{q^{\text{in}},\delta_{k,k_i(\mathbf{c})}} \sum_{j \in \partial_i} \delta_{k',k_j(\mathbf{c})} \delta_{q^{\text{out}},\delta_{k',k_i(\mathbf{c})}} \sum_{j \in \partial_i} \delta_{k,k_j(\mathbf{c})} \\ &= \frac{1}{N} \sum_i \left[\delta_{k,k_i(\mathbf{c})} \delta_{q^{\text{in}},\sum_{j \in \partial_i} \delta_{k',k_j(\mathbf{c})}} + (1 - \delta_{k,k_i(\mathbf{c})}) \delta_{q^{\text{in}},0} \right] \\ &\quad \times \left[\delta_{k',k_i(\mathbf{c})} \delta_{q^{\text{out}},\sum_{j \in \partial_i} \delta_{k,k_j(\mathbf{c})}} + (1 - \delta_{k',k_i(\mathbf{c})}) \delta_{q^{\text{out}},0} \right] \\ &= \sum_n p(n|\mathbf{c}) \left[\delta_{k,k(n)} \delta_{q^{\text{in}},\sum_{s \leq k(n)} \delta_{k',\xi^s(n)}} + (1 - \delta_{k,k(n)}) \delta_{q^{\text{in}},0} \right] \\ &\quad \times \left[\delta_{k',k(n)} \delta_{q^{\text{out}},\sum_{s \leq k(n)} \delta_{k,\xi^s(n)}} + (1 - \delta_{k',k(n)}) \delta_{q^{\text{out}},0} \right]. \quad (15) \end{aligned}$$

The two marginals of (15) are

$$p_{\text{in}}^{kk'}(q) = \sum_n p(n|\mathbf{c}) \left[\delta_{k,k(n)} \delta_{q,\sum_{s \leq k(n)} \delta_{k',\xi^s(n)}} + (1 - \delta_{k,k(n)}) \delta_{q,0} \right], \quad (16)$$

$$p_{\text{out}}^{kk'}(q) = \sum_n p(n|\mathbf{c}) \left[\delta_{k',k(n)} \delta_{q,\sum_{s \leq k(n)} \delta_{k,\xi^s(n)}} + (1 - \delta_{k',k(n)}) \delta_{q,0} \right]. \quad (17)$$

Hence $p_{\text{in}}^{kk'}(q) = p_{\text{out}}^{k'k}(q)$, as expected. The average degree $\bar{q}^{kk'} = \sum_{q^{\text{in}},q^{\text{out}}} q^{\text{in}} p^{kk'}(q^{\text{in}}, q^{\text{out}}) = \sum_{q^{\text{in}},q^{\text{out}}} q^{\text{out}} p^{kk'}(q^{\text{in}}, q^{\text{out}})$ of the graph $\beta^{kk'}$ can be written, using identity (10) and the symmetry of $W(k, k'|\mathbf{c})$, as

$$\bar{q}^{kk'} = \sum_n p(n|\mathbf{c}) \delta_{k(n),k} \sum_{s \leq k(n)} \delta_{k',\xi^s(n)} = \bar{k}(\mathbf{c}) W(k, k'|\mathbf{c}). \quad (18)$$

If $k = k'$, the decomposition matrix $\beta^{kk'}$ is symmetric. Here we find

$$k_i(\beta^{kk}) = \delta_{k,k_i(\mathbf{c})} \sum_{j \in \partial_i} \delta_{k,k_j(\mathbf{c})}. \quad (19)$$

Hence the degree distribution of β^{kk} becomes

$$\begin{aligned} p^{kk}(q) &= \frac{1}{N} \sum_i \delta_{q, \delta_{k, k(i)}} \sum_{j \in \partial_i} \delta_{k, k(j)} \\ &= \frac{1}{N} \sum_i \left[\delta_{k, k(i)} \delta_{q, \sum_{j \in \partial_i} \delta_{k, k(j)}} + (1 - \delta_{k, k(i)}) \delta_{q, 0} \right] \\ &= \sum_n p(n|\mathbf{c}) \left[\delta_{k, k(n)} \delta_{q, \sum_{s \leq k(n)} \delta_{k, \xi^s(n)}} + (1 - \delta_{k, k(n)}) \delta_{q, 0} \right]. \end{aligned} \quad (20)$$

The average degree in β^{kk} is therefore

$$\bar{q}^{kk} = \sum_n p(n|\mathbf{c}) \delta_{k(n), k} \sum_{s \leq k(n)} \delta_{k, \xi^s(n)} = \bar{k}(\mathbf{c}) W(k, k|\mathbf{c}). \quad (21)$$

4. Entropy of ensembles of bipartite graphs

Here we calculate the leading orders in N of the entropy per node (6) for ensembles of bipartite graphs with prescribed (and possibly distinct) degree distributions in the two node sets. This is not only a novel result in itself, but will also form the seed of the entropy calculation for ensembles with constrained neighbourhoods in a subsequent section.

In a bipartite ensemble the N nodes can be divided into two disjoint sets $A, B \subseteq \{1, \dots, N\}$ such that $c_{ij} = 0$ if $i, j \in A$ or $i, j \in B$, leaving only links *between* A and B . This constraint implies that there is a mapping from the set of bipartite graphs on $\{1, \dots, N\}$ to the set of directed graphs on $\{1, \dots, N\}$, defined by assigning to each bipartite link the direction of flow from A to B . This allows us to draw upon results on directed graphs derived in [2]. The directed graph \mathbf{c}' associated with the bipartite graph \mathbf{c} would have

$$j \in B \text{ or } i \in A: \quad c'_{ij} = 0, \quad (22)$$

$$j \in A \text{ and } i \in B: \quad c'_{ij} = c_{ij} \quad (23)$$

and hence the in- and out-degree sequence $\vec{k} = ((k_1^{\text{in}}, k_1^{\text{out}}), \dots, (k_N^{\text{in}}, k_N^{\text{out}}))$ of \mathbf{c}' can be expressed in terms of the degree sequence \mathbf{k} of \mathbf{c} via

$$i \in A: \quad \vec{k}_i = (k_i^{\text{in}}, k_i^{\text{out}}) = (0, k_i), \quad (24)$$

$$i \in B: \quad \vec{k}_i = (k_i^{\text{in}}, k_i^{\text{out}}) = (k_i, 0). \quad (25)$$

This mapping can be shown to be bijective between the ensemble of all bipartite graphs with a given degree sequence, and the ensemble of all directed graphs with the appropriately chosen directed degree sequence. The directed graph will have the joint degree distribution

$$p(q^{\text{in}}, q^{\text{out}}) = \frac{|A|}{N} \delta_{q^{\text{in}}, 0} p_A(q^{\text{out}}) + \left(1 - \frac{|A|}{N}\right) p_B(q^{\text{in}}) \delta_{q^{\text{out}}, 0} \quad (26)$$

with the degree distributions $p_A(k) = |A|^{-1} \sum_{i \in A} \delta_{k, k(i)}$ and $p_B(k) = |B|^{-1} \sum_{i \in B} \delta_{k, k(i)}$ in the sets A and B of the bipartite graph. Our bipartite ensemble is one in which we describe the distributions $p_A(k)$ and $p_B(k)$, together with the probability $f \in [0, 1]$ for a node to be in subset A , and we forbid links within the sets A or B . Conservation of links demands that the

two distributions cannot be independent, but must obey $\bar{q} = (1-f) \sum_q q p_B(q) = f \sum_q q p_A(q)$, where \bar{q} is the average degree. Hence, the entropy of any degree-constrained bipartite ensemble can be calculated by application of (6), (7) to the associated ensemble of degree-constrained directed graphs, with $X_i = (\tau_i, k_i)$. Here $\tau_i \in \{A, B\}$ gives the subset assignment of a node. We then find

$$S = \sum_{\tau, k} \left[\prod_i p(\tau_i, k_i) \right] S(\tau, k) - f \log f - (1-f) \log(1-f) - f \sum_k p_A(k) \log p_A(k) - (1-f) \sum_k p_B(k) \log p_B(k) \quad (27)$$

with

$$p(\tau, k) = f \delta_{\tau, A} p_A(k) + (1-f) \delta_{\tau, B} p_B(k), \quad (28)$$

$$S(\tau, k) = \frac{1}{N} \log \sum_{\mathbf{c}} \left(\prod_{i, \tau_i=A} \delta_{\bar{k}_i, (0, k_i)} \right) \left(\prod_{i, \tau_i=B} \delta_{\bar{k}_i, (k_i, 0)} \right). \quad (29)$$

The latter quantity follows from the calculation in [7], with the short-hand $\pi_{\bar{q}}(q) = e^{-\bar{q}q}/q!$ and modulo terms that vanish for $N \rightarrow \infty$:

$$\begin{aligned} S(\tau, k) &= \bar{q} \left[\log(N/\bar{q}) + 1 \right] + \sum_q \left[f \delta_{q,0} + (1-f) p_B(q) \right] \log \pi_{\bar{q}}(q) \\ &\quad + \sum_q \left[f p_A(q) + (1-f) \delta_{q,0} \right] \log \pi_{\bar{q}}(q) \\ &= \bar{q} \log(N/\bar{q}) + f \sum_q p_A(q) \log \pi_{\bar{q}}(q) + (1-f) \sum_q p_B(q) \log \pi_{\bar{q}}(q). \end{aligned} \quad (30)$$

This then leads to our final result for the entropy per node of tailored bipartite graph ensembles, with imposed bipartite degree distributions $p_A(k)$ and $p_B(k)$, average degree \bar{k} , and a fraction f of nodes in the set A (modulo vanishing orders in N):

$$\begin{aligned} S &= \bar{k} \log(N/\bar{k}) - f \log f - (1-f) \log(1-f) \\ &\quad - f \sum_k p_A(k) \log \left(\frac{p_A(k)}{\pi_{\bar{k}}(k)} \right) - (1-f) \sum_k p_B(k) \log \left(\frac{p_B(k)}{\pi_{\bar{k}}(k)} \right). \end{aligned} \quad (31)$$

If the sets A and B were to be specified explicitly (as opposed to only their relative sizes), the contribution $S_f = -f \log f - (1-f) \log(1-f)$ would disappear from the above formula.

5. Entropy of ensembles with constrained neighbourhoods

We now turn to the Shannon entropy per node (6) of the ensemble (2) in which for the observables $X_i(\mathbf{c})$ we choose the local neighbourhood $n_i(\mathbf{c})$ defined in (1). For this we need to calculate the leading orders of $S(\mathbf{n}) = N^{-1} \log \sum_{\mathbf{c}} \delta_{\mathbf{n}, \mathbf{n}(\mathbf{c})}$. We now use the one-to-one relationship between a graph \mathbf{c} and its decomposition $\mathbf{c} = \sum_{qq'} \beta^{qq'}$, to write

$$S(\mathbf{n}) = \frac{1}{N} \log \sum_{\{\beta^{kk'}\}} \delta_{\mathbf{n}, \mathbf{n}(\mathbf{c})}. \quad (32)$$

The next argument is the key to our ability to evaluate the entropy. It involves translating the constraint $\mathbf{n} = \mathbf{n}(\mathbf{c})$ into constraints on the decomposition matrices $\beta^{kk'}$. Let us define the sets of nodes in \mathbf{c} which have the same degree, viz. $I_k(\mathbf{n}) = \{i \leq N | k_i(\mathbf{c}) = k\}$. The constraint $\mathbf{n} = \mathbf{n}(\mathbf{c})$ in (32) prescribes:

- (i) all the sets I_k of nodes with a given degree
- (ii) for each node $i \in I_k$ which sets $I_{k'}$ this node is (possibly multiply) connected to.

Hence the constraint $\mathbf{n} = \mathbf{n}(\mathbf{c})$ specifies exactly the in- and out-degree sequences of all decomposition matrices $\beta^{kk'}$ of \mathbf{c} , which we will denote as $\vec{q}^{kk'} = (q^{\text{in},kk'}, q^{\text{out},kk'})$, and whose distributions we have already calculated in (15) and (20). We thus see that (32) can be written as

$$S(\mathbf{n}) = \frac{1}{N} \log \sum_{\{\beta^{kk'}\}} \prod_{kk'} \delta_{\vec{q}_n^{kk'}, \vec{q}(\beta^{kk'})}, \quad (33)$$

in which $\vec{q}_n^{kk'}$ are the in- and out-degree sequences that are imposed by the local environment sequence \mathbf{n} on the decomposition matrix $\beta^{kk'}$, and whose distributions are known to be (15), (20). Using the symmetry $(\beta^{kk'})^\dagger = \beta^{k'k}$ we may now write

$$\begin{aligned} S(\mathbf{n}) &= \frac{1}{N} \log \left[\left(\prod_{k < k'} \sum_{\beta^{kk'}} \delta_{\vec{q}_n^{kk'}, \vec{q}(\beta^{kk'})} \right) \left(\prod_k \sum_{\beta^{kk}} \delta_{\vec{q}_n^{kk}, \vec{q}(\beta^{kk})} \right) \right] \\ &= \sum_{k < k'} \left\{ \frac{1}{N} \log \sum_{\beta^{kk'}} \delta_{\vec{q}_n^{kk'}, \vec{q}(\beta^{kk'})} \right\} + \sum_k \left\{ \frac{1}{N} \log \sum_{\beta^{kk}} \delta_{\vec{q}_n^{kk}, \vec{q}(\beta^{kk})} \right\}. \end{aligned} \quad (34)$$

We see that the entropy $S(\mathbf{n})$ can be written as the sum of the entropies of sub-ensembles, which are the decomposition matrices $\beta^{kk'}$ with prescribed degree sequences. The second sum in (34) is over nondirected ensembles, the first over directed ones. The sub-entropies were all calculated, respectively, in [1] and [2]⁴. The entropy of an N -node nondirected random graph ensemble with degree sequence \mathbf{q} was found to be (modulo terms that vanish for $N \rightarrow \infty$):

$$S_{\mathbf{q}} = \frac{1}{N} \log \sum_{\mathbf{c}} \delta_{\mathbf{q}, \mathbf{q}(\mathbf{c})} = \frac{1}{2} \bar{q} \left[\log(N/\bar{q}) + 1 \right] + \sum_q p(q) \log \pi_{\bar{q}}(q) \quad (35)$$

in which $\bar{q} = N^{-1} \sum_i q_i$ and $\pi_{\bar{q}}(q)$ is the Poisson distribution with average \bar{q} . The entropy of an N -node directed random graph ensemble with in- and out-degree sequence \vec{q} was found to be (modulo terms that vanish for $N \rightarrow \infty$):

$$\begin{aligned} S_{\vec{q}} &= \frac{1}{N} \log \sum_{\mathbf{c}} \delta_{\vec{q}, \vec{q}(\mathbf{c})} \\ &= \bar{q} \left[\log(N/\bar{q}) + 1 \right] + \sum_{q^{\text{in}}, q^{\text{out}}} p(q^{\text{in}}, q^{\text{out}}) \log \left[\pi_{\bar{q}}(q^{\text{in}}) \pi_{\bar{q}}(q^{\text{out}}) \right]. \end{aligned} \quad (36)$$

The above entropies depend in leading orders only on the degree distributions (as opposed to the degree sequences), and since these distributions were already calculated (15) and (20), we can simply insert (36) and (36) into (34), with the correct distributions (15) and (20), and find

⁴ In [1, 2] the entropies were carried out for ensembles with prescribed degree distributions, but it was shown that, in analogy with (6), this is simply the sum of the Shannon entropy of the degree distributions and the entropy of the corresponding ensemble with prescribed sequences.

an expression that depends only on the local environment distribution $p(n) = N^{-1} \sum_i \delta_{n,n_i}$:

$$\begin{aligned}
S(\mathbf{n}) &= \sum_{k < k'} \left\{ \bar{q}^{kk'} \left[\log \left(N / \bar{q}^{kk'} \right) + 1 \right] + \sum_{q^{\text{in}}, q^{\text{out}}} p^{kk'}(q^{\text{in}}, q^{\text{out}}) \log \left[\pi_{\bar{q}^{kk'}}(q^{\text{in}}) \pi_{\bar{q}^{kk'}}(q^{\text{out}}) \right] \right\} \\
&\quad + \sum_k \left\{ \frac{1}{2} \bar{q}^{kk} \left[\log \left(N / \bar{q}^{kk} \right) + 1 \right] + \sum_q p^{kk}(q) \log \pi_{\bar{q}^{kk}}(q) \right\} \\
&= \frac{1}{2} \sum_{k \neq k'} \left\{ \bar{k} W(k, k') \left[\log \left(N / \bar{k} W(k, k') \right) - 1 \right] \right. \\
&\quad \left. + 2 \bar{k} W(k, k') \log \left[\bar{k} W(k, k') \right] - \sum_q \left[p_{\text{in}}^{kk'}(q) + p_{\text{out}}^{kk'}(q) \right] \log q! \right\} \\
&\quad + \frac{1}{2} \sum_k \left\{ \bar{k} W(k, k) \left[\log \left(N / \bar{k} W(k, k) \right) - 1 \right] \right. \\
&\quad \left. + 2 \bar{k} W(k, k) \log \left[\bar{k} W(k, k) \right] - 2 \sum_q p^{kk}(q) \log q! \right\} \\
&= \frac{1}{2} \bar{k} \left[\log \left(N / \bar{k} \right) - 1 \right] + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
&\quad - \sum_q \left[\frac{1}{2} \sum_{k \neq k'} \left[p_{\text{in}}^{kk'}(q) + p_{\text{out}}^{kk'}(q) \right] + \sum_k p^{kk}(q) \right] \log q! \\
&= \frac{1}{2} \bar{k} \left[\log \left(N / \bar{k} \right) - 1 \right] + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
&\quad - \sum_q \sum_n p(n) \sum_{k, k'} \left[\delta_{k, k(n)} \delta_{q, \sum_{s \leq k(n)} \delta_{k', \xi^s(n)}} + (1 - \delta_{k', k(n)}) \delta_{q, 0} \right] \log q! \\
&= \frac{1}{2} \bar{k} \left[\log \left(N / \bar{k} \right) - 1 \right] + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
&\quad - \sum_n p(n) \sum_k \log \left[\left(\sum_{s \leq k(n)} \delta_{k, \xi^s(n)} \right)! \right]. \tag{37}
\end{aligned}$$

Insertion of this result into the general formula (6) gives us an analytical expression for the Shannon entropy of the random graph ensemble with prescribed distribution $p(n)$ of local neighbourhoods, modulo terms that vanish for $N \rightarrow \infty$. This expression is fully explicit, since \bar{k} and $W(k, k')$ are both determined by the distribution $p(n)$, via $\bar{k} = \sum_n p(n) k(n)$ and (10) respectively:

$$\begin{aligned}
S &= \frac{1}{2} \bar{k} \left[\log \left(N / \bar{k} \right) - 1 \right] + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
&\quad - \sum_n p(n) \log p(n) - \sum_n p(n) \sum_k \log \left[\left(\sum_{s \leq k(n)} \delta_{k, \xi^s(n)} \right)! \right]. \tag{38}
\end{aligned}$$

6. Entropy of ensembles of networks with specified generalized degree distribution

In this section we consider an ensemble of nondirected networks with a specified generalized degree distribution $p(k, m) = N^{-1} \sum_i \delta_{k, k_i(c)} \delta_{m, m_i(c)}$, where $k_i(c) = \sum_j c_{ij}$ and $m_i = \sum_{jk} c_{ij} c_{jk}$. Previous work [3] began this calculation, and reached (in leading order) the intermediate form set out below:

$$S = \frac{1}{2} \bar{k} \left[\log(N/\bar{k}) + 1 \right] - \sum_{k,m} p(k, m) \log \left(\frac{p(k, m)}{\pi_{\bar{k}}(k)} \right) + \sum_{k,m} p(k, m) \log \left(\sum_{\xi^1, \dots, \xi^k} \delta_{m, \sum_{s=1}^k \xi^s} \prod_{s=1}^k \gamma(k, \xi^s) \right) \quad (39)$$

\bar{k} indicates the average degree; $\pi_{\bar{k}}(k)$ is the Poissonian distribution with average degree \bar{k} . The sum inside the logarithm in the final term of (39) runs over all sets of k nonnegative integers $\xi^1 \dots \xi^k$. The function $\gamma(\dots)$ is defined as the nonnegative solution to the following self-consistency relation:

$$\gamma(k, k') = \sum_{m'} \frac{k'}{\bar{k}} p(k', m') \left[\frac{\sum_{\xi^1, \dots, \xi^{k'-1}} \delta_{m'-k, \sum_{s=1}^{k'-1} \xi^s} \prod_{s=1}^{k'-1} \gamma(k', \xi^s)}{\sum_{\xi^1, \dots, \xi^{k'}} \delta_{m', \sum_{s=1}^{k'} \xi^s} \prod_{s=1}^{k'} \gamma(k', \xi^s)} \right]. \quad (40)$$

This equation does not yield to a straightforward solution, and can only be evaluated numerically or in certain special cases. Without a physical interpretation of $\gamma(k, k')$, this intermediate answer is limited in how much insight it can provide. We will now show how the entropy can be expressed in terms of measurable quantities.

Within the notion of taking node properties from a specified distribution there is the interesting question of whether it is even possible to realize a network with the given local topological properties. For simple degrees, this question is generally considered to not be material in the $N \rightarrow \infty$ limit. For finite N , the Erdos–Gallai theorem famously gives a necessary and sufficient condition for a degree sequence to be graphical. This question is not explicitly considered within the statistical mechanics approach—which makes it particularly interesting that there arises the self-consistency equation (40), which seems to have a natural interpretation from the point of view of graphicality.

Our strategy is to derive an expression for the (observable) degree–degree correlations $W(k, k')$, and show that these can be expressed in terms of the order parameter $\gamma(k, k')$ that appears in equation (39). We calculate the average of this quantity in our tailored ensembles of the form (2), where we now define topological characteristics by specifying a generalized degree distribution $p(k, m)$. We follow closely the steps taken in [3], and write for the ensemble a specified generalized degree sequence (\mathbf{k}, \mathbf{m}) :

$$\begin{aligned}
 & W(k, k') \\
 &= \frac{1}{N\bar{k}} \sum_{\mathbf{c}} P(\mathbf{c}|\mathbf{k}, \mathbf{m}) \sum_{rs} c_{rs} \delta_{k, \sum_t c_{rt}} \delta_{k', \sum_t c_{st}} = \frac{1}{N^2} \sum_{rs} \delta_{k, k_r} \delta_{k', k_s} \\
 &\quad \int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) - i(\theta_r + \theta_s + \phi_r k_s + \phi_s k_r)} \prod_{i < j} \left[1 + \frac{\bar{k}}{N} \left(e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} - 1 \right) \right] \\
 &\quad \times \frac{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m})} \prod_{i < j} \left[1 + \frac{\bar{k}}{N} \left(e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} - 1 \right) \right]}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m})} \prod_{i < j} \left[1 + \frac{\bar{k}}{N} \left(e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} - 1 \right) \right]} \\
 &\quad + \mathcal{O}\left(\frac{1}{N}\right) \\
 &= \frac{1}{N^2} \sum_{rs} \delta_{k, k_r} \delta_{k', k_s} \frac{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) - i(\theta_r + \theta_s + \phi_r k' + \phi_s k)} + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} + \dots}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m})} + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} + \dots} \\
 &\quad + \mathcal{O}\left(\frac{1}{N}\right) \\
 &= \frac{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} \left(\frac{1}{N^2} \sum_{rs} \delta_{k, k_r} \delta_{k', k_s} e^{-i(\theta_r + \theta_s + \phi_r k' + \phi_s k)} \right)}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}}} + \mathcal{O}\left(\frac{1}{N}\right) \\
 &= \frac{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)} \left(\frac{1}{N} \sum_r \delta_{k, k_r} e^{-i(\theta_r + \phi_r k')} \right) \left(\frac{1}{N} \sum_s \delta_{k', k_s} e^{-i(\theta_s + \phi_s k)} \right)}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot \mathbf{k} + \phi \cdot \mathbf{m}) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}}} \\
 &\quad + \mathcal{O}\left(\frac{1}{N}\right) \\
 &= \frac{\int \{dP d\hat{P}\} e^{N\mathcal{Y}[P, \hat{P}]} \left(\int d\theta d\phi P(\theta, \phi, k) e^{-i\theta - i\phi k'} \right) \left(\int d\theta d\phi P(\theta, \phi, k') e^{-i\theta - i\phi k} \right)}{\int \{dP d\hat{P}\} e^{N\mathcal{Y}[P, \hat{P}]}} \\
 &\quad + \mathcal{O}\left(\frac{1}{N}\right). \tag{41}
 \end{aligned}$$

Taking the limit $N \rightarrow \infty$ therefore gives

$$\begin{aligned}
 \lim_{N \rightarrow \infty} W(k, k') &= \left(\int d\theta d\phi P(\theta, \phi, k) e^{-i\theta - i\phi k'} \right) \\
 &\quad \times \left(\int d\theta d\phi P(\theta, \phi, k') e^{-i\theta - i\phi k} \right) \Big|_{\text{saddle-point } \{P, \hat{P}\} \text{ of } \mathcal{Y}} \tag{42}
 \end{aligned}$$

in which the function $\mathcal{Y}[P, \hat{P}]$ is identical to that found in [3]. Using the formulae in [3] that relate to the definition of the order parameter $\gamma(k, k')$, we then obtain for $N \rightarrow \infty$ the unexpected simple but welcome relation

$$W(k, k') = \gamma(k, k')\gamma(k', k). \quad (43)$$

A similar, although slightly more involved, calculation leads to an expression for the joint distribution $W(k, m; k', m')$; see the appendix for details.

Our final aim is to use identity (43) to resolve equation (39) into observable quantities. Consider the nontrivial term in (39)

$$\Gamma = \sum_{k,m} p(k, m) \log \left(\sum_{\xi^1, \dots, \xi^k} \prod_{s=1}^k \gamma(k, \xi^s) \delta_{m, \sum_{s=1}^k \xi^s} \right). \quad (44)$$

At this point of the calculation, the effect of factorising across nodes has been to break the expression down into terms which, for every generalized degree (k, m) , enumerate all the possible ways of dividing m second neighbours between k first neighbours. The term inside the logarithm sums for each k over all configurations $\{\xi^1 \dots \xi^k\}$ which meet the condition $\sum_{s=1}^k \xi^s = m$. To formalize this idea, we may re-aggregate the expression for any graphically realizable distribution $p(k, m)$ to write

$$\begin{aligned} \Gamma &= \frac{1}{N} \log \left\{ \prod_{k,m} \left(\sum_{\xi^1, \dots, \xi^k} \prod_{s=1}^k \gamma(k, \xi^s) \delta_{m, \sum_{s=1}^k \xi^s} \right)^{Np(k,m)} \right\} \\ &= \frac{1}{N} \log \prod_i \left(\sum_{\xi_i^1, \dots, \xi_i^{k_i}} \left[\prod_{s=1}^{k_i} \gamma(k_i, \xi_i^s) \right] \delta_{m_i, \sum_{s=1}^{k_i} \xi_i^s} \right) \\ &= \frac{1}{N} \log \left\{ \sum_{\xi_1^1, \dots, \xi_1^{k_1}} \dots \sum_{\xi_N^1, \dots, \xi_N^{k_N}} \left(\prod_i \delta_{m_i, \sum_{s=1}^{k_i} \xi_i^s} \right) \prod_i \prod_{s=1}^{k_i} \gamma(k_i, \xi_i^s) \right\}. \quad (45) \end{aligned}$$

We can now see that the separate terms precisely enumerate all the permutations of degrees and neighbour-degrees for networks with a generalized degree sequence consistent with any pair (k, m) appearing $Np(k, m)$ times. The Kronecker deltas $\delta_{m_i, \sum_{s=1}^{k_i} \xi_i^s}$ tell us that each ξ_i^s in any nonzero term is to be interpreted as the degree of a node $j \in \partial_i$, and must therefore appear also as the left argument in another factor of the type $\gamma(k_j, \cdot)$. This insight allows the expression to be substantially simplified, since we already know that $\gamma(k, k')\gamma(k', k) = W(k', k)$, where $W(k', k)$ is the correlation between degrees of connected nodes. Hence, any nonvanishing contribution to the sum over all neighbourhoods inside the logarithm of (45) will be equal to a repeated product of factors $W(k, k')$, with different (k, k') . Since we also know that the number of links between nodes with degree combination (k, k') equals $N\bar{k}W(k, k')$ in leading order in N , we conjecture that in leading order we may make the following replacement inside (45):

$$\prod_i \prod_{s=1}^{k_i} \gamma(k_i, \xi_i^s) \rightarrow \prod_{k,k'} W(k, k')^{\frac{\bar{k}N}{2} W(k,k')}, \quad (46)$$

(where the factor $\frac{1}{2}$ in the exponent reflects the fact that two $\gamma(\cdot, \cdot)$ factors combine to form each factor $W(\cdot, \cdot)$.) With this conjecture we obtain, in leading order in N :

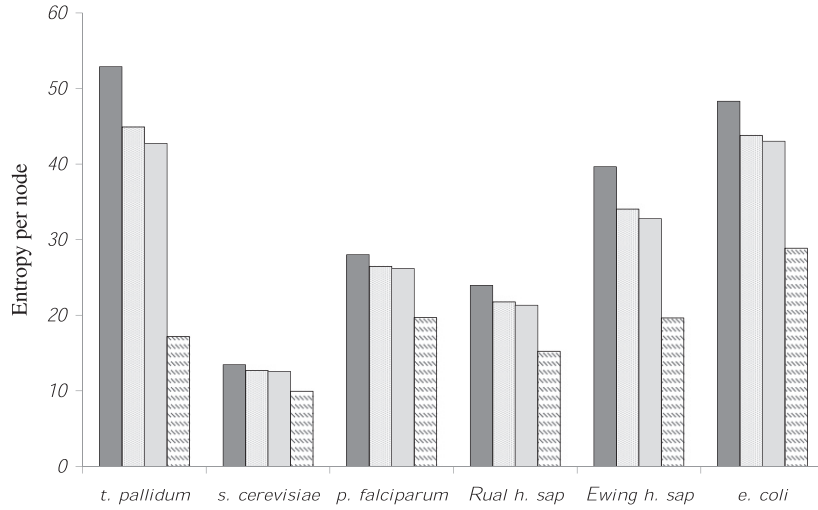


Figure 2. Results from applying equation (48) for the entropy of random graph ensembles—where the constraints are taken to match the relevant topological observables of networks from [8–13]. From left to right the bars correspond to entropy per node of random graph ensembles tailored to match: average degree, degree distribution, degree–degree correlation and generalized degrees. The generalized degrees constraint can be seen to be onerous for several of these networks. This observation is based on the large reduction in entropy between the generalized degree constrained ensemble, and any of the other ensembles considered.

$$\begin{aligned}
 \Gamma &= \frac{1}{N} \log \left\{ \sum_{\xi_1^1 \dots \xi_1^{k_1}} \dots \sum_{\xi_N^1 \dots \xi_N^{k_N}} \left(\prod_i \delta_{m_i, \sum_{s=1}^{k_i} \xi_i^s} \right) \right\} + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
 &= \sum_{k, m} p(k, m) \log \left(\sum_{\xi^1, \dots, \xi^k} \delta_{m, \sum_{s=1}^k \xi^s} \right) + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k'). \quad (47)
 \end{aligned}$$

This implies that (39) simplifies to

$$\begin{aligned}
 S &= \frac{1}{2} \bar{k} \left[\log(N/\bar{k}) + 1 \right] + \frac{1}{2} \bar{k} \sum_{k, k'} W(k, k') \log W(k, k') \\
 &\quad - \sum_{k, m} p(k, m) \log \left(\frac{p(k, m)}{\pi_{\bar{k}}(k)} \right) + \sum_{k, m} p(k, m) \log \left(\sum_{\xi^1, \dots, \xi^k} \delta_{m, \sum_{s=1}^k \xi^s} \right). \quad (48)
 \end{aligned}$$

Figure 2 applies equation (48), and related results from [1], to real protein–protein interaction datasets, in order to demonstrate how these expressions can quantify the degree to which each topological property constrains the ensemble. In appendix B we look at some simple synthetic examples, where (48) can be compared with direct enumeration, or where the self-consistency equation (40) is soluble.

7. Conclusion

Ensembles of tailored random graphs are extremely useful constructions in the modelling of complex interacting particle systems in biology, physics, computer science, economics and the social sciences. They allow us to quantify topological features of such systems and reason quantitatively about their complexity, as well as define and generate useful random proxies for realistic networks in statistical mechanical analyses of processes.

In this paper we have derived, in leading two orders in N , explicit expressions for the Shannon entropies of different types of tailored random graph ensembles, for which no such expressions had yet been obtained. This work builds on and extends the ideas and techniques developed in the three papers [1–3], which use path integral representations to achieve link factorization in the various summations over graphs. We show in this paper how the new ensemble entropies can often be calculated by efficient use and combination of earlier results.

The first class of graph ensembles we studied consists of bipartite nondirected graphs with prescribed (and possibly nonidentical) distributions of degrees for the two node subsets. This case is handled by a bijective mapping from the ensemble of bipartite graphs with a specified degree distribution to the ensemble of directed graphs with the associated specified directed degree distribution, for which formulae are available. The second class consists of graphs with prescribed distributions of local neighbourhoods, where the neighbourhood of a node is defined as its own degree plus the values of the degrees of its immediate neighbours. This problem was solved using a decomposition in terms of bipartite graphs, building on the previous result. The final class of graphs, for which the entropy had in the past only partially been calculated, consist of graphs with prescribed distributions of generalized degrees, i.e. of ordinary degrees plus the total number of length-two paths starting in the specified nodes. Here we derive two novel and exact identities linking the order parameters to macroscopic observables, which lead to an explicit entropy formula based on a plausible but not yet proven conjecture,

Since completing this work, our attention has been drawn to a preprint [14] which considers the question of the entropy of random graph ensembles constrained with a given distribution of neighbourhoods by a probability theory route, via an adapted configuration model. In that case, the neighbourhoods were specified as graphlets of an arbitrary depth. [14] also retrieves the entropy of an ensemble constrained with a specified degree distribution, as originally derived by [1].

Acknowledgements

ESR gratefully acknowledges financial support from the Biotechnology and Biological Sciences Research Council of the United Kingdom under grant no BB/H018409/1.

Appendix A. Generalized degree correlation kernel for ensembles with prescribed generalized degrees

The generalized quantity $W(k, m; k', m')$ in the ensemble with prescribed generalized degree distributions $p(k, m)$ can be calculated along the same lines as the calculation of $W(k, k')$ in the main text. It is defined as

$$W(k, m; k', m' | \mathbf{c}) = \frac{1}{N\bar{k}} \sum_{ij} c_{ij} \delta_{k, \sum_{\ell} c_{i\ell}} \delta_{k', \sum_{\ell} c_{j\ell}} \delta_{m, \sum_{\ell} c_{i\ell} k_{\ell}} \delta_{m', \sum_{\ell} c_{j\ell} k'_{\ell}} \quad (\text{A.1})$$

and its ensemble average takes the form

$$\begin{aligned} & W(k, m; k', m') \\ &= \frac{\left[\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot k + \phi \cdot m) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}} \left(\frac{1}{N^2} \sum_{rs} \delta_{k, k_r} \delta_{k', k_s} \delta_{m, m_r} \delta_{m', m_s} e^{-i(\theta_r + \theta_s + \phi_r k' + \phi_s k)} \right) \right]}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot k + \phi \cdot m) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}} + \mathcal{O}\left(\frac{1}{N}\right)} \\ &= \frac{\left[\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot k + \phi \cdot m) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}} \left(\frac{1}{N} \sum_r \delta_{k, k_r} \delta_{m, m_r} e^{-i(\theta_r + \phi_r k')} \right) \left(\frac{1}{N} \sum_s \delta_{k', k_s} \delta_{m', m_s} e^{-i(\theta_s + \phi_s k)} \right) \right]}{\int_{\pi}^{\pi} d\theta d\phi e^{i(\theta \cdot k + \phi \cdot m) + \frac{\bar{k}}{2N} \sum_{ij} e^{-i(\theta_i + \theta_j + \phi_i k_j + \phi_j k_i)}} + \mathcal{O}\left(\frac{1}{N}\right)}. \end{aligned} \quad (\text{A.2})$$

Now we will want to introduce a generalized order parameter, namely

$$P(\theta, \phi, k, m) = \frac{1}{N} \sum_r \delta_{k, k_r} \delta_{m, m_r} \delta(\theta - \theta_r) \delta(\phi - \phi_r), \quad (\text{A.3})$$

The previous order parameter used in the calculation of $W(k, k')$ is a marginal of this, via $P(\theta, \phi, k) = \sum_m P(\theta, \phi, k, m)$. This definition will give us

$$\begin{aligned} W(k, m; k', m') &= \left(\int_{-\pi}^{\pi} d\theta d\phi P(\theta, \phi, k, m) e^{-i\theta - i\phi k'} \right) \\ &\quad \times \left(\int_{-\pi}^{\pi} d\theta d\phi P(\theta, \phi, k', m') e^{-i\theta - i\phi k} \right), \end{aligned} \quad (\text{A.4})$$

$$W(k, k') = \sum_{mm'} W(k, m; k', m') \quad (\text{A.5})$$

in which the new order parameter and its conjugate are to be solved by extremization of the generalized surface

$$\begin{aligned} \Psi[P, \hat{P}] &= i \sum_{km} \int_{-\pi}^{\pi} d\theta d\phi \hat{P}(\theta, \phi, k, m) P(\theta, \phi, k, m) \\ &\quad + \sum_{km} P(k, m) \log \int_{-\pi}^{\pi} d\theta d\phi e^{i(\theta k + \phi m - \hat{P}(\theta, \phi, k, m))} \\ &\quad + \frac{1}{2} \bar{k} \int d\theta d\phi d\theta' d\phi' \sum_{kk'mm'} P(\theta, \phi, k, m) P(\theta', \phi', k', m') e^{-i(\theta + \theta' + \phi k' + \phi' k)}. \end{aligned} \quad (\text{A.6})$$

Variation of Ψ gives the following saddle-point equations

$$\hat{P}(\theta, \phi, k, m) = i\bar{k}e^{-i\theta} \int d\theta' d\phi' \sum_{k'm'} P(\theta', \phi', k', m') e^{-i(\theta'+\phi k'+\phi'k)}, \quad (\text{A.7})$$

$$P(\theta, \phi, k, m) = P(k, m) \frac{e^{i(\theta k + \phi m - \hat{P}(\theta, \phi, k, m))}}{\int_{-\pi}^{\pi} d\theta' d\phi' e^{i(\theta' k + \phi' m - \hat{P}(\theta', \phi', k, m))}}. \quad (\text{A.8})$$

Clearly $\hat{P}(\theta, \phi, k, m) = \hat{P}(\theta, \phi, k)$ (i.e. it is independent of m). We may therefore substitute $\hat{P}(\theta, \phi, k) = i\bar{k}e^{-i\theta} \hat{P}(\phi, k)$ and find

$$\hat{P}(\phi, k) = \int d\theta' d\phi' \sum_{k'm'} P(\theta', \phi', k', m') e^{-i(\theta'+\phi k'+\phi'k)}, \quad (\text{A.9})$$

$$P(\theta, \phi, k, m) = P(k, m) \frac{e^{i(\theta k + \phi m) + \bar{k}e^{-i\theta} \hat{P}(\phi, k)}}{\int_{-\pi}^{\pi} d\theta' d\phi' e^{i(\theta' k + \phi' m) + \bar{k}e^{-i\theta'} \hat{P}(\phi', k)}}. \quad (\text{A.10})$$

We observe as before in [3] that

$$\begin{aligned} \int_{-\pi}^{\pi} d\theta P(\theta, \phi, k) e^{-i\theta} &= \sum_m P(k, m) \frac{\int_{-\pi}^{\pi} d\theta e^{i(\theta(k-1) + \phi m) + \bar{k}e^{-i\theta} \hat{P}(\phi, k)}}{\int_{-\pi}^{\pi} d\theta d\phi' e^{i(\theta k + \phi' m) + \bar{k}e^{-i\theta} \hat{P}(\phi', k)}} \\ &= \sum_m P(k, m) \frac{\sum_{\ell \geq 0} \frac{\bar{k}^{\ell} \hat{P}^{\ell}(\phi, k)}{\ell!} \int_{-\pi}^{\pi} d\theta e^{i(\theta(k-1-\ell) + \phi m)}}{\sum_{\ell \geq 0} \frac{\bar{k}^{\ell} \hat{P}^{\ell}(\phi', k)}{\ell!} \int_{-\pi}^{\pi} d\theta d\phi' e^{i(\theta k + \phi' m - \ell\theta)}} \\ &= \sum_m P(k, m) \frac{\frac{\bar{k}^{k-1} \hat{P}^{k-1}(\phi, k)}{(k-1)! e^{i\phi m}}}{\frac{\bar{k}^k \hat{P}^k(\phi', k)}{k!} \int_{-\pi}^{\pi} d\phi' e^{i\phi' m}} \\ &= \sum_m \frac{k}{\bar{k}} P(k, m) \frac{\hat{P}^{k-1}(\phi, k) e^{i\phi m}}{\int_{-\pi}^{\pi} d\phi' \hat{P}^k(\phi', k) e^{i\phi' m}}. \end{aligned} \quad (\text{A.11})$$

Hence

$$\hat{P}(\phi, k) = \sum_{k'm'} \frac{k'}{\bar{k}} P(k', m') e^{-i\phi k'} \frac{\int_{-\pi}^{\pi} d\phi' \hat{P}^{k'-1}(\phi', k') e^{i\phi'(m'-k)}}{\int_{-\pi}^{\pi} d\phi' \hat{P}^{k'}(\phi', k') e^{i\phi' m'}}. \quad (\text{A.12})$$

After writing $\hat{P}(\phi, k) = \sum_{k'} \gamma(k, k') e^{-i\phi k'}$ we recover our familiar equation

$$\gamma(k, k') \gamma(k', k) = \frac{k'}{\bar{k}} \sum_m P(k', m) \frac{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m, \sum_{n \leq k'} k_n} \delta_{kk_n}}{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m, \sum_{n \leq k'} k_n}}. \quad (\text{A.13})$$

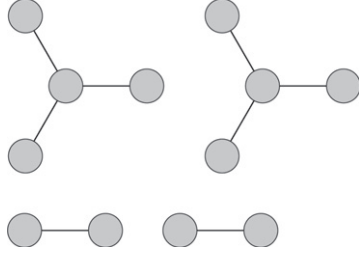


Figure B1. A validation example with nonconstant degree-degree correlation.

But now we can also work out the generalized kernel:

$$\begin{aligned}
W(k, m; k', m') &= \left(\int_{-\pi}^{\pi} d\theta d\phi P(\theta, \phi, k, m) e^{-i\theta - i\phi k'} \right) \left(\int_{-\pi}^{\pi} d\theta d\phi P(\theta, \phi, k', m') e^{-i\theta - i\phi k} \right) \\
&= \frac{kk'}{\bar{k}^2} P(k, m) P(k', m') \left(\frac{\int_{-\pi}^{\pi} d\phi \hat{P}^{k-1}(\phi, k) e^{i\phi(m-k)}}{\int_{-\pi}^{\pi} d\phi \hat{P}^k(\phi, k) e^{i\phi m}} \right) \\
&\quad \times \left(\frac{\int_{-\pi}^{\pi} d\phi \hat{P}^{k'-1}(\phi, k') e^{i\phi(m'-k)}}{\int_{-\pi}^{\pi} d\phi \hat{P}^{k'}(\phi, k') e^{i\phi m'}} \right) \\
&= \frac{kk'}{\bar{k}^2} \frac{P(k, m) P(k', m')}{\gamma(k, k') \gamma(k', k)} \left(\frac{\sum_{k_1 \dots k_k} \left[\prod_{n=1}^k \gamma(k, k_n) \right] \delta_{m, \sum_{n \leq k} k_n} \delta_{k, k'}}{\sum_{k_1 \dots k_k} \left[\prod_{n=1}^k \gamma(k, k_n) \right] \delta_{m, \sum_{n \leq k} k_n}} \right) \\
&\quad \times \left(\frac{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m', \sum_{n \leq k'} k_n} \delta_{k', k}}{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m', \sum_{n \leq k'} k_n}} \right). \tag{A.14}
\end{aligned}$$

We know that $W(k, k') = \gamma(k, k') \gamma(k', k)$, and that $P(k, m) k / \bar{k} = W(k, m)$, so this can be simplified to

$$\begin{aligned}
W(k, m; k', m') &= \frac{W(k, m) W(k', m')}{W(k, k')} \left(\frac{\sum_{k_1 \dots k_k} \left[\prod_{n=1}^k \gamma(k, k_n) \right] \delta_{m, \sum_{n \leq k} k_n} \delta_{k, k'}}{\sum_{k_1 \dots k_k} \left[\prod_{n=1}^k \gamma(k, k_n) \right] \delta_{m, \sum_{n \leq k} k_n}} \right) \\
&\quad \times \left(\frac{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m', \sum_{n \leq k'} k_n} \delta_{k', k}}{\sum_{k_1 \dots k_{k'}} \left[\prod_{n=1}^{k'} \gamma(k', k_n) \right] \delta_{m', \sum_{n \leq k'} k_n}} \right). \tag{A.15}
\end{aligned}$$

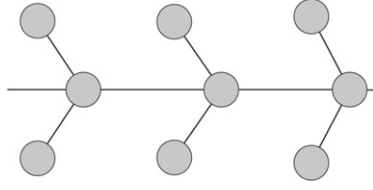


Figure B2. A validation example based on a connected network.

Appendix B. Some synthetic examples evaluated directly and with the derived formulae

Below we analyse some synthetic networks, where it is possible to evaluate the entropy directly, and compare this to the results of applying equation (48). Consider a network based on figure B1 as a repeating motif.

The self-consistency relations defined in equation (40) are straightforward

$$\gamma^2(1, 1) = \frac{p(1, 1)}{\bar{k}} = \frac{1}{4} \quad \gamma(1, 3)\gamma(3, 1) = \frac{3p(3, 3)}{\bar{k}} = \frac{3}{8} \quad (\text{B.1})$$

from which it follows that Γ defined in equation (44) is

$$\Gamma = \frac{\log \gamma(1, 1)}{3} + \frac{\log \gamma(1, 3)}{2} + \frac{\log \gamma^3(3, 1)}{6} = \frac{1}{6} \log \left[\frac{1}{4} \left(\frac{3}{8} \right)^3 \right]. \quad (\text{B.2})$$

This evaluates to the same as the expression for Γ proposed in equation (47)

$$\begin{aligned} \Gamma &= \sum_{k,m} p(k, m) \log \left(\sum_{\xi^1, \dots, \xi^k} \delta_{m, \sum_{s=1}^k \xi^s} \right) + \frac{1}{2} \bar{k} \sum_{k,k'} W(k, k') \log W(k, k') \\ &= \frac{8}{6} \frac{1}{2} \left[\frac{1}{4} \log \frac{1}{4} + \frac{3}{8} \log \frac{3}{8} + \frac{3}{8} \log \frac{3}{8} \right] + 0 = \frac{1}{6} \log \left[\frac{1}{4} \left(\frac{3}{8} \right)^3 \right], \end{aligned} \quad (\text{B.3})$$

where $W(k, k')$ is defined in equation (5). A combinatorial argument can be used to count the number of networks in the ensemble—which is the same as evaluating the partition function $Z(\cdot)$ that was defined in equation (2). The number of configurations can be obtained by counting the number of labellings of the diagram, and then dividing by the symmetries, which evaluates to:

$$Z(\mathbf{k}, \mathbf{m}) = \frac{\frac{N}{3}! \frac{N}{2}!}{\frac{N}{6}! 2!^{N/6} 3!^{N/6}}.$$

Applying Stirling's formula, it follows that in leading order

$$\begin{aligned} \frac{1}{N} \log Z(\mathbf{k}, \mathbf{m}) &= \frac{1}{3} \log \left(\frac{N}{3} \right) + \frac{1}{2} \log \left(\frac{N}{2} \right) - \frac{1}{6} \log \left(\frac{N}{6} \right) - \frac{1}{6} \log(12) - \frac{4}{6} \\ &= \frac{4}{6} \log N - \frac{1}{3} [\log 3 + 2 \log 2] - \frac{2}{3}. \end{aligned} \quad (\text{B.4})$$

This matches the analytic result, evaluated with the help of equation (48). The link between (48) and $Z(\mathbf{k})$ is defined in equations (6) and (7). Consider the ensemble of networks defined by the general degree sequence set out in figure B2. By combinatorics, argue that there are $\frac{N}{3}!$

orderings of the centre nodes, and $\frac{2N}{3}!$ orderings of the remaining nodes, divided by $2^{\frac{N}{3}}$ for symmetry. It follows that in leading order

$$\begin{aligned} \frac{1}{N} \log Z(\mathbf{k}, \mathbf{m}) &= \frac{1}{N} \log \left(\frac{\frac{N}{3}! \frac{2N}{3}!}{2^{\frac{N}{3}}} \right) \\ &= \frac{1}{3} \log \frac{N}{3} + \frac{2}{3} \log \frac{2N}{3} - \frac{1}{3} \log 2 - 1 \\ &= \log N - \log 3 + \frac{1}{3} \log 2 - 1. \end{aligned} \quad (\text{B.5})$$

The analytic formula evaluates to

$$\begin{aligned} \frac{1}{N} \log Z(\mathbf{k}, \mathbf{m}) &= \log N + \log 2 - 1 - \frac{2}{3} \log 2 - \log 3 \\ &= \log N + \frac{1}{3} \log 2 - \log 3 - 1. \end{aligned} \quad (\text{B.6})$$

References

- [1] Annibale A, Coolen A C, Fernandes L P, Fraternali F and Kleinjung J 2009 *J. Phys. A: Math. Theor.* **42** 485001
- [2] Roberts E S, Schlitt T and Coolen A C 2011 *J. Phys. A: Math. Theor.* **44** 5002
- [3] Bianconi G, Coolen A C and Vicente C J P 2008 *Phys. Rev. E* **78** 016114
- [4] Faudree R J, Gould R J and Lesniak L M 1992 *Discrete Appl. Math.* **37-38** 179–91
- [5] Rogers T, Vicente C P, Takeda K and Castillo I P 2010 *J. Phys. A: Math. Theor.* **43** 195002
- [6] Newman M 2009 *Networks: an Introduction* (Oxford: Oxford University Press)
- [7] Roberts E S and Coolen A C 2012 *Phys. Rev. E* **85** 046103
- [8] Titz B, Rajagopala S V, Goll J, Häuser R, McKevitt M T, Palzkill T and Uetz P 2008 *PLoS One* **3** e2292
- [9] Ito T, Chiba T, Ozawa R, Yoshida M, Hattori M and Sakaki Y 2001 *Proc. Natl. Acad. Sci. USA* **98** 4569–74
- [10] LaCount D J *et al* 2005 *Nature* **438** 103–7
- [11] Rual J F *et al* 2005 *Nature* **437** 1173–8
- [12] Ewing R M *et al* 2007 *Mol. Syst. Biol.* **3** 89
- [13] Arifuzzaman M *et al* 2006 *Genome Res.* **16** 686–91
- [14] Bordenave C and Caputo P 2013 arXiv:1308.5725