

Random Graph Models with Hidden Color *

BO SÖDERBERG

Complex Systems Div., Dept. for Theoretical Physics, Lund University,
Sölvegatan 14 A, S-223 62 LUND, Sweden

LU TP 03-34

Submitted to Acta Physica Polonica B

We demonstrate how to generalize two of the most well-known random graph models, the *classic random graph*, and *random graphs with a given degree distribution*, by the introduction of hidden variables in the form of extra degrees of freedom, color, applied to vertices or stubs (half-edges). The color is assumed unobservable, but is allowed to affect edge probabilities. This serves as a convenient method to define very general classes of models within a common unifying formalism, and allows for a non-trivial edge correlation structure.

PACS numbers: 02.50.-r, 64.60.-i, 89.75.Fb

1. Introduction

The availability of data on real-world networks, *e.g.* from information technology and molecular biology, has seen a dramatic increase in the last decades. This has led to a correspondingly increased interest in the theoretical modelling of networks.

Typically the growth of a real-world network is not entirely deterministic but contains stochastic elements, and statistical models are required that are conveniently formulated in terms of *ensembles* of graphs. Typical real-world networks are not static but change with time, and much of the focus has been on *dynamical* models, where one attempts to describe the growth and evolution of a network.

Here we will focus on *static* random graph models, describing a snapshot of a network in terms of a fixed ensemble of graphs, without regard to how the network was formed. By a *random graph* we will mean a member of such

* Presented at the Workshop on Random Geometry in Krakow, May 2003

an ensemble. In particular, we will be mostly interested in *sparse* random graphs, where the typical vertex degree does not grow with the size of the graph.

There is a vast spectrum of such models around. Some of these are not entirely random, in the sense that they are based on an underlying regular network – i.e. a lattice – which is then modified in a random fashion.

Our focus will be on *purely random* graphs, where such an underlying regularity is absent. A number of more or less unrelated models of this type have been investigated, and it would obviously be desirable to devise a unified description in terms of a *general class of ensembles*, where more specialized models appear as special cases of one and the same general formalism, while maintaining the computability of local and global graph characteristics of interest, such as degree distributions, small subgraph abundancies, component size distributions, and global connectivity properties.

The most well-known purely random model is the *classic Random Graph* of Erdős and Rényi [1], to be referred to as *RG*. In its sparse version it is defined as follows. For a given set of N nodes, every pair of nodes is connected by an edge independently with probability $p = c/N$ in terms of a given parameter c that asymptotically defines the average degree. This model has many interesting properties, such as an asymptotically Poissonian degree distribution, and a phase transition at $c = 1$, above which a giant connected component is formed. However, it fails to describe most real-world networks.

A more general model that has been much studied is *Random Graphs with a given Degree Distribution* [2, 3, 4, 5], or Degree-driven Random Graphs (*DRG*), where an asymptotic degree distribution is given, suitably transformed into a definite degree sequence for a given graph size. In terms of this a random graph is defined as a uniformly random member drawn from the set of graphs having the given degree sequence, possibly subject to additional constraints (*e.g.* by demanding the graph to be simple, *i.e.* non-degenerate). DRG models suffer from an intrinsic lack of edge correlations, atypical of real-world networks; as a result, they are often referred to as uncorrelated random graphs.

In a sequence of papers [6, 7, 8], I have explored the use of hidden coloring, either of vertices or of stubs, to define more general random graph models. The resulting models can be seen as colored extensions of RG and DRG. As will be shown in this paper, the hidden color provides a convenient means for defining very general classes of random graph models, where much of the limitations of the uncolored models can be done away with, while the computability of interesting properties is maintained.

The considered classes of models will be compared with respect to a few basic properties: The degree distribution, the abundancy of arbitrarily

given small subgraphs, the size distribution of connected components, and the phase transition where a giant component appears.

The plan of the paper is as follows. In section 2, we will discuss a few fundamental concepts needed in the subsequent sections. In section 3, we will review the definitions of the models to be considered. A comparative analysis of a selected set of characteristics, as derived in the different model classes, is presented in section 4. Section 5, finally, contains a concluding discussion.

2. Basic Concepts and Methods

All models to be considered in this paper will be of sparse random graphs, where the degrees (connectivities) of vertices stay finite as the size N of the graph grows to infinity. In particular, this means that the total number of edges will scale as N , and that the probability of a connection between an arbitrary pair of nodes will scale as $1/N$.

A simple local characteristic of a graph ensemble is its *degree distribution*, $\{p_m\}$. This is often conveniently described in terms of its *generating function*,

$$H(x) \equiv \sum_m p_m x^m. \quad (1)$$

It obviously satisfies $H(1) = 1$, and yields upon repeated differentiation at $x = 1$ successive *combinatorial moments* of the degree,

$$H'(1) = \langle m \rangle, \quad H''(1) = \langle m(m-1) \rangle, \quad H'''(1) = \langle m(m-1)(m-2) \rangle, \quad (2)$$

etc., while the individual p_m can be obtained by repeated differentiation at $x = 0$.

Generating functions of this type are convenient when analyzing a probability distribution P_k of an integer variable k that is the *sum* of several independent contributions, $k = \sum_i k_i$, in which case the generating function $f(z) = \sum_k P_k z^k$ for the distribution of k is simply the *product* of the corresponding generating functions for the distribution of each contribution.

3. The Models

Here follows a brief introduction to the models to be considered.

3.1. The Classic Model – RG

The classic random graph (RG) [1] has been thoroughly analyzed over the years [9, 10]. It is a model of simple (non-degenerate) labelled graphs with a given set of N nodes, although it can be easily extended to include

also non-simple graphs [11]. We will consider its sparse version. It comes in two essentially equivalent versions, one with a fixed number of edges, the other with a fixed probability for each possible edge; we will stick to the latter.

Sparse RG has a single real parameter $c > 0$ controlling the abundance of edges. For a given graph size N and a given value of c , an ensemble of graphs is defined as follows. Each of the $N(N - 1)/2$ pairs of distinct nodes independently is connected by an undirected edge with a common probability $p = c/N$ (assuming $N > c$).

3.2. Inhomogeneous Random Graphs – IRG

The classic RG model as described above can be generalized in a straightforward way by assigning color to vertices and allowing edge probabilities to be color-sensitive; the resulting class of models will be referred to as **IRG**, for inhomogeneous random graphs [6].

A definite IRG model is specified in terms of

- a color space, taken as $[1, \dots, K]$;
- a *color distribution* $\{r_a > 0, a = 1, \dots, K\}$, with $\sum_a r_a = 1$;
- a real, symmetric *color preference matrix* $\mathbf{c} = \{c_{ab} \geq 0\}$.

For a given graph size N , such a model is implemented as follows.

1. Assign to each node independently a random color a , drawn from the given distribution $\{r_a\}$.
2. Connect each pair of distinct nodes independently with probability c_{ab}/N , where a and b are the respective colors of the two nodes.

By considering the color as unobservable or hidden, the resulting ensemble of colored graphs yields a specific ensemble of plain graphs, distinct from an RG ensemble, as will be shown below. The role of the hidden color is to enable non-trivial edge correlations. IRG defines a class of graph ensembles much more general than RG.

3.3. Random graphs with a given degree distribution – DRG

The classic RG model is limited to a Poissonian degree distribution. A more general class of models that has recently attracted the attention of several workers in the statistical physics community is random graphs with a

given degree distribution [2, 4, 5], to be referred to as *DRG* (for degree-driven random graphs).¹ This approach allows for an arbitrary degree distribution.

There are two common variants of DRG. One is given by restricting the ensemble to simple (non-degenerate) graphs, where self-couplings and multiple connections are banned. In the other (the configuration model) one allows for degenerate graphs. For ease of analysis, we will focus on the latter version where degeneracies are allowed.

Some notation: A node with degree m is considered to possess m *stubs*, each of which defines a point of attachment for an edge endpoint. The total number of stubs $M = \sum_i m_i$ in a graph obviously must be even, being equal to the total number of edge endpoints, i.e. twice the number of edges.

A specific DRG model is defined by specifying an arbitrary degree distribution $\{p_m\}$. For a fixed graph size N , the corresponding ensemble is implemented as follows:²

1. For each node, draw its degree independently from the given distribution. Redo until the total sum of the degrees is even.
2. Define random edges by performing a completely random pairing within the resulting even-numbered set of M stubs.

This leads to an ensemble of pseudographs, where degeneracies may appear, in the form of self-connections (tadpoles) or multiple edges between the same pair of nodes.

The random stub pairing is reminiscent of the combinatorics associated with Gaussian integrals; indeed, a relation exists between DRG models and certain miniature field theories [12, 13].

3.4. DRG plus color – CDRG

Also the DRG class of models can be generalized, by utilizing a coloring of *stubs*, which turns out to be the most natural choice, and then allowing the stub pairing to be color-sensitive [7, 8]. The resulting very general class of models will be referred to as CDRG, for Colored DRG.

With colored stubs, it is natural to consider the color-specific stub content of a node, its *colored degree*. With K colors to choose between, the colored degree is conveniently represented by an integer vector $\mathbf{m} = (m_1, \dots, m_K)$, with the individual elements m_a counting the number of stubs with color a . Obviously the plain degree m is obtained by summing

¹ Variants of this approach have been referred under various names, such as *equilibrium random graphs* and *uncorrelated random graphs*.

² We disregard impossible cases of an odd N with a degree distribution supporting only odd degrees.

up the elements of the colored degree, $m = \sum_a m_a = \mathbf{1} \cdot \mathbf{m}$, in terms of the uniform vector $\mathbf{1} = (1, \dots, 1)$.

Then it is also natural to consider the probability distribution of such colored degrees, a *colored degree distribution* $\{p_{\mathbf{m}}\}$. Such a distribution can be represented by a *multivariate generating function*, $H(\mathbf{x}) = \sum_{\mathbf{m}} p_{\mathbf{m}} \mathbf{x}^{\mathbf{m}}$, where $\mathbf{x}^{\mathbf{m}} = \prod_a x_a^{m_a}$, satisfying the normalizing condition $H(\mathbf{1}) = 1$. From this, *multivariate combinatorial moments* can be derived by repeated differentiation at $\mathbf{x} = \mathbf{1}$, e.g. $\partial_a H(\mathbf{x} = \mathbf{1}) = \langle m_a \rangle$, $\partial_a \partial_b H(\mathbf{x} = \mathbf{1}) = \langle m_a m_b - m_a \delta_{ab} \rangle$, etc. A specific CDRG model is defined by specifying

- A color space, taken as $[1, \dots, K]$;
- A colored degree distribution $\{p_{\mathbf{m}}\}$;
- A real, symmetric color preference matrix $\mathbf{T} = \{T_{ab} \geq 0\}$, such that $\mathbf{T} \langle \mathbf{m} \rangle = \mathbf{1}$.

We will for simplicity assume that the colored degree distribution is such that all moments are defined.

For a given graph size N , such a model is implemented as follows.

1. For each node, draw its colored degree independently from the given distribution. Redo until the total sum of the (plain) degrees is even.
2. Define random edges by performing a weighted random pairing within the resulting even-numbered set of M stubs, such that the probability for each of the $(M - 1)!!$ possible pairings has a statistical weight proportional to the product over all edges of a factor given by T_{ab} , where a, b are the colors of the stubs it connects.

This class of models obviously collapses to DRG for the case of a single color, in which case the matrix \mathbf{T} collapses to a single number, given by $1/\langle m \rangle$ by virtue of the constraint $\mathbf{T} \langle \mathbf{m} \rangle = \mathbf{1}$.

The constraint on \mathbf{T} is convenient for the forthcoming analysis, and ensures that the total number of ab -edges asymptotically approaches the value $N \langle m_a \rangle T_{ab} \langle m_b \rangle$, which upon summing over b yields the correct asymptotic number of a -stubs as $N \langle m_a \rangle$.

The combinatorics of the weighted random pairing yields the following asymptotic results. The probability that two arbitrary stubs with known colors a, b will be paired with each other is T_{ab}/N . This implies that the probability for two random nodes with respective colored degrees \mathbf{m}, \mathbf{m}' will be connected is $\sum_{ab} m_a T_{ab} m'_b / N$, which reduces (as it must!) to $\langle m \rangle / N$ if the degrees are not known.

A less general class of models, similar in spirit to CDRG but restricted to homogeneously colored vertices (i.e. with vertex coloring rather than stub coloring), has also been investigated [14].

4. Analysis

Below follows a comparative analysis, where we review a selected set of local and global characteristics for random graphs as drawn from models of the different types, with a focus on the asymptotic limit $N \rightarrow \infty$. For a more detailed analysis, we refer to the paper [8] and references therein.

4.1. Asymptotic Degree Distributions

First we will derive the resulting asymptotic degree distributions, where not defined in the model specifications.

4.1.1. RG degree distributions

In a graph drawn from an RG model as described above, each node has $N - 1$ possible connections, each independently realized with probability c/N . Thus, the degree m of a random node will obey a binomial distribution, $\binom{N-1}{m} (c/N)^m (1 - c/N)^{N-1-m}$. As $N \rightarrow \infty$ with fixed c , this approaches an asymptotic distribution $\{p_m\}$, given by a *Poissonian* with average c ,

$$p_m = e^{-c} \frac{c^m}{m!}, \quad (3)$$

with the corresponding generating function $H(x) = e^{c(x-1)}$.

4.1.2. IRG degree distributions

Choose a random node in a large graph from an IRG ensemble. It has the color a with probability r_a . For large N there are $\sim Nr_b$ other nodes with color b ; each of these is connected to the chosen node independently with probability c_{ab}/N . Thus, for a node of given color a , we asymptotically expect its number of b -neighbors to follow a Poissonian distribution with average $c_{ab}r_b$, and its total degree to follow a Poissonian with average $C_a \equiv \sum_b c_{ab}r_b$. Averaging over a yields the asymptotic degree distribution

$$p_m = \sum_a r_a e^{-C_a} \frac{C_a^m}{m!}, \quad (4)$$

with the generating function

$$H(x) = \sum_a r_a e^{C_a(x-1)}, \quad (5)$$

which describes a *Poissonian mix*. This implies the following convexity constraint on the possible degree distributions:

$$p_m^2 \leq \frac{m+1}{m} p_{m-1} p_{m+1}, \quad (6)$$

for $m > 0$. Conversely, any degree distribution obeying this constraint can, at least in principle, be realized with a suitable IRG model, possibly with infinitely many colors.

4.1.3. DRG degree distributions

The degree distribution is considered given in a DRG model, and so is in principle free to choose. For ease of analysis, we shall restrict our considerations to cases where all moments $\langle m^n \rangle$ exist, barring power-behaved distributions, which otherwise are interesting in their own right.

4.1.4. CDRG degree distributions

In a CDRG model, a colored degree distribution is given, from which the plain degree distribution can be extracted directly. Its generating function $H(x)$ is obtained simply by evaluating the multivariate generating function (with the same name) for the colored degree distribution with a homogeneous argument, $H(x) = H(x\mathbf{1}) \equiv H(x, \dots, x)$.

Since the colored degree distribution is free to choose, so is the plain one, and there are obviously many CDRG models with a given degree distribution.

4.2. Small Subgraph Statistics

The combinatorial moments of the degree distribution are simply related to the expected numbers of subgraphs in the form av stars. More general local characteristics can be expressed in terms of the number of copies of an arbitrary small graph γ found as subgraphs of a large random graph G . We will be interested in the expected number of copies in the asymptotic limit $N \rightarrow \infty$.

The clustering properties of a graph are often analyzed in terms of the probability of two neighbors of a node to be connected; this is seen to be related to the number of simple triangles, i.e. the number of subgraphs γ in the form of a mutually connected triple of nodes.

Thus, assume an arbitrary small connected graph γ to be given, having $v \ll N$ vertices and $e \ll N$ edges. We can estimate its expected number $\langle n_\gamma \rangle$ of distinct occurrences as a subgraph in a random graph G of size N as follows.

A particular possible embedding of γ in G is defined by mapping the ordered set of v nodes in γ onto a target set given by an ordered v -subset of the N vertices in G . There are $N!/(N-v)! \approx N^v$ such sets. However, for a target set to define a valid subgraph position, each edge in γ must be mapped onto an existing edge in the target set.

For the models considered, the expected count $\langle n_\gamma \rangle$ can be derived from Feynman-like rules, with model-specific vertex and node factors as well as the usual symmetry factors.

4.2.1. Small subgraphs in RG

The RG model describes simple random graphs, which can have only simple subgraphs. For each of the $\sim N^v$ possible embeddings of a simple γ , the probability for the corresponding set of e target edges to exist is $(c/N)^e$.

Thus, naively, the expected number of occurrences should be $N^{v-e}c^e$. If γ has a non-trivial isomorphism group, i.e. a symmetry under some permutation of its vertices, the naive result has to be divided by the order S_γ of the symmetry group. This leaves us with the following simple rules for the asymptotically expected subgraph count $\langle n_\gamma \rangle$.

- For each node in γ , associate a factor N .
- For each edge in γ , associate a factor c/N .
- Multiply the node and edge factors, and divide the result by the symmetry factor S_γ .

Since γ is assumed connected, we have $e \geq v - 1$, and $e - v + 1$ counts its number of loops. Thus, the expected count scales as $O(N)$ for a *tree*, and as $O(1)$ for a one-loop γ , while it vanishes asymptotically for γ with several loops. This is typical of a sparse random graph – loops are scarce.

As an example illustrating the lack of correlations in an RG ensemble, consider subgraphs in the form of a v -chain, i.e. a set of v nodes connected in an open chain; the expected counts show a simple geometric behaviour, $\langle n_\gamma \rangle = Nc^{v-1}/2$.

4.2.2. Small subgraphs in IRG

Also in IRG, graphs are simple, so also here, we must assume γ to be simple. Generalizing the arguments used for RG, we get the following rules for the asymptotically expected number $\langle n_\gamma \rangle$.

- Associate with each node in γ an independent color a , and a corresponding factor Nr_a .
- Associate with each edge in γ a factor c_{ab}/N , where a, b are the node colors at its endpoints.
- Multiply all node and edge factors, sum over the node colors, and divide the result by the symmetry factor S_γ .

Again, expected counts for tree subgraphs scale as $O(N)$, and those for connected one-loop subgraphs as $O(1)$.

Non-trivial edge correlations are possible in IRG, as illustrated by v -chain subgraphs, where the hidden color in an IRG ensemble enables the expected counts to deviate from the simple geometric behavior found for a plain RG ensemble; instead it takes the form of a mix of geometric sequences.

4.2.3. Small subgraphs in DRG

Since a DRG ensemble of the kind we are considering allows for degeneracies, we will have to consider also possibly degenerate subgraphs, with loops of length one or two. Since subgraphs with loops are suppressed due to the sparsity, just as in RG and IRG, degeneracies will turn out not to be very important.

The expected number of copies of γ with a fixed set of target nodes can be calculated as follows. Consider a node in the target set with actual degree m , that defines the target for a node with degree k in γ . The corresponding k target edges can be chosen among the m existing ones in $m_k \equiv m!/(m-k)!$ distinct ways. This can be shown to yield the following rules for calculating the asymptotic $\langle n_\gamma \rangle$ for a DRG model.

- Associate with each node with k stubs in γ a factor $N \langle m_k \rangle$.
- Associate with each edge in γ a factor $1/(N \langle m \rangle)$.
- Multiply the node and edge factors, and divide the result by the symmetry factor S_γ , including possible contributions from edge permutations and flips for the case of a non-simple γ .

Here, $\langle m_k \rangle$ stands for the k th combinatorial moment, defined by $\partial_z^k H(z=1) = \langle m(m-1)\dots(m-k+1) \rangle$ (see eq. (2)).

For a v -chain, the expected count becomes $N \langle m \rangle^{3-v} \langle m(m-1) \rangle^{v-2} / 2$, displaying simple geometric behavior just as for the case of RG, illustrating the lack of edge correlations in DRG.

4.2.4. Small subgraphs in CDRG

To analyze the subgraph statistics for a CDRG model, we will need the colored generalizations of the combinatorial moments,

$$E_{abc\dots} \equiv \partial_a \partial_b \partial_c \dots H(\mathbf{z} = \mathbf{1}), \quad (7)$$

with the lowest ones given by $E_a \equiv \langle m_a \rangle$, $E_{ab} \equiv \langle m_a m_b - m_a \delta_{ab} \rangle$, etc. Generalizing the argument used for DRG, one can derive the following rules [8] for the asymptotically expected subgraph counts in a CDRG model.

- Associate with each stub in γ an independent color label.
- Associate with each node in γ having k stubs a factor $NE_{abc\dots}$, where $a, b, c \dots$ are the k associated color labels.
- Associate with each edge in γ a factor T_{ab}/N , where a, b are the color labels associated with the stubs at its endpoints.
- Multiply the node and edge factors, sum over the associated colors, and divide the result by the symmetry factor S_γ , including possible contributions from edge permutations and flips for the case of a non-simple γ .

Note how these reduce to the DRG rules for the case of a single color.

For the case of a v -chain, the expected count becomes a mix of geometric sequences, showing how the coloring also for CDRG enables non-trivial edge correlations, just as was the case for IRG.

4.3. Connected Component Sizes

Next we turn to an analysis of the global connectivity characteristics of a random graph. These are simplest described in terms of the sizes of the connected components of the graph.

Thus, we will be interested in the size distribution P_n of a connected component of a random graph, as revealed from a randomly chosen initial node by recursively exploring edges leading to new nodes until the entire component is revealed. For sparse random graphs in the asymptotic limit $N \rightarrow \infty$, any finite component is almost surely a tree, since any extra connections will be suppressed by factors of $1/N$. Thus, the revelation of such a component can be described as a *branching process*, with properties depending on the specific model considered.

The asymptotic component size distribution is conveniently analyzed in terms of its generating function,

$$g(z) \equiv \sum_n P_n z^n. \quad (8)$$

For the models considered, $g(z)$ or a set of related functions will satisfy recursive equations that determine the sought distribution.

4.3.1. Component sizes in RG

For an RG model, the component size distribution can be estimated as follows, as long as the component remains small. For each revealed node i , the number k of branches to new nodes obeys a Poissonian distribution

asymptotically, $p_k \sim e^{-c}c^k/k!$, since there are $\sim N$ remaining unrevealed nodes, each of which connects to i with probability c/N .

This yields the recursive equation $g(z) = ze^{-c} \sum_k c^k g(z)^k/k!$, to be understood as follows. The initial factor of z accounts for the initial node, while each term in the sum describes the case where it has a distinct number k of neighbors. The factor $e^{-c}c^k/k!$ represents the probability for this case, and the factor $g(z)^k$ encodes the fact that each of the k neighbors defines a subtree statistically identical to the full tree.

The recursion can be simplified to read

$$g(z) = ze^{c(g(z)-1)}, \quad (9)$$

which should be interpreted as an *iterated map* for the value of g for a given value of z , a *stable fixed point* of which defines the physical value.

As a curiosity, eq. (9) can be written as $F(cg(z)) = zF(c)$, with $F(c) = ce^{-c}$, with the explicit solution $g(z) = F^{-1}(zF(c))/c$, with the inverse of F defined from the restriction $F(c), |c| \leq 1$. Taylor-expanding the inverse yields the exact solution $P_n = \frac{(nce^{-c})^n}{cnn!}$ for $n \geq 1$ for the asymptotic component size distribution, with the large- n behaviour $P_n \rightarrow \frac{(ce^{1-c})^n}{\sqrt{2\pi cn^{3/2}}}$, decaying exponentially for $c \neq 1 \Rightarrow ce^{1-c} < 1$, but only as a power for $c = 1$.

4.3.2. Component sizes in IRG

For an IRG model, the asymptotic component size distribution, and thus its generating function $g(z)$, will obviously be an average over the result conditional upon a particular color of the initial node, and we can write

$$g(z) = \sum_a r_a g_a(z), \quad (10)$$

where $g_a(z)$ is conditional upon initial color a ; these satisfy the following recursive relations.

$$g_a(z) = z \exp \left[\sum_b c_{ab} r_b (g_b(z) - 1) \right], \quad (11)$$

an obvious generalization of the corresponding RG result, eq. (9). This cannot in general be solved exactly, but can be analyzed using numerical and/or series expansion methods.

4.3.3. Component sizes in DRG

Next we wish to obtain the asymptotic size distribution $\{P_n\}$ in a DRG model. As before, the sparsity forces finite components to take the form of tree.

The generating function $H(x)$ for the degree distribution will turn out to be convenient. Of interest is also the degree distribution of a node reached by following a random edge. This yields a weighting of nodes by their degree, resulting in the modified distribution $q_m = mp_m/\langle m \rangle$. The generating function for its *remaining* degree (disregarding the incoming stub) becomes $H'(x)/H'(1)$.

With $g(z)$ as before the generating function for the size distribution of the entire component, let $h(z)$ be the analogous generating function for the size distribution of a subtree found by following an edge. Then $g(z)$ can be expressed in terms of $h(z)$ as

$$g(z) = z \sum_m p_m h(z)^m \equiv zH(h(z)), \tag{12}$$

to be interpreted as follows. The explicit factor of z represents the first node. It has m outgoing edges with probability p_m , each of which represents a subtree and yields a factor $h(z)$; see fig. 1 for a graphical illustration. By

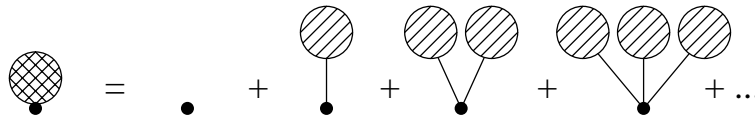


Fig. 1. $g(z)$ in terms of $h(z)$ for a DRG model, illustrating eq. (12).

a similar argument, $h(z)$ satisfies the recursion

$$h(z) = z \sum_m \frac{mp_m}{\langle m \rangle} h(z)^{m-1} = z \frac{H'(h(z))}{H'(1)}, \tag{13}$$

as depicted in fig. 2. Note that for the case of a Poissonian degree distri-

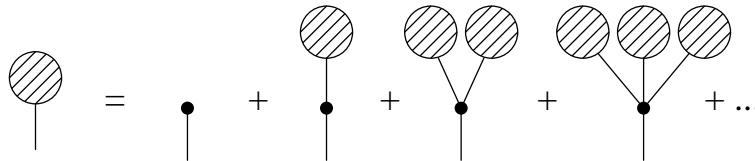


Fig. 2. Illustration of the recursive relation (13) for $h(z)$.

bution, the recursion simplifies to the RG result. Indeed, the Poissonian restriction of DRG is asymptotically equivalent to a version of RG allowing for non-simple graphs.

4.3.4. Component sizes in CDRG

Finally, we wish to obtain the asymptotic size distribution $\{P_n\}$, and its associated generating function $g(z)$, in a CDRG model.

Here, the multivariate generating function $H(\mathbf{x})$ for the colored degree distribution will be needed. Of interest is also the colored degree distribution $q_{\mathbf{m}|a}$ of a node reached by following a random edge emanating from a stub of given color a . This is given by $q_{\mathbf{m}|a} = \sum_b T_{ab} m_b p_{\mathbf{m}}$. It follows that the generating function for the distribution of its *remaining* colored degree (where the incoming stub is neglected) is $\sum_b T_{ab} \partial_b H(\mathbf{x})$.

With $g(z)$ having its usual meaning, we will denote by $h_a(z)$ the analogous generating function for the size distribution of a subtree found by following an edge emanating from a stub of color a . Then, generalizing eq. (12), $g(z)$ can be expressed in terms of $\mathbf{h}(z) = (h_1(z), \dots, h_K(z))$ as

$$g(z) = z \sum_{\mathbf{m}} p_{\mathbf{m}} \prod_a h_a(z)^{m_a} \equiv zH(\mathbf{h}(z)), \quad (14)$$

with the following interpretation. The explicit factor of z accounts for the first node, which has a colored degree \mathbf{m} with probability $p_{\mathbf{m}}$; each stub of color a represents a subtree and yields a factor $h_a(z)$. The argument generalizes the one used for DRG, as depicted in fig. 1.

By a similar argument, $\mathbf{h}(z)$ satisfies the coupled recursion

$$h_a(z) = z \sum_b T_{ab} \partial_b H(\mathbf{h}(z)), \quad (15)$$

generalizing the DRG relation, eq. (13), depicted in fig. 2.

4.4. The Appearance of the Giant

For all models, $g(z)$ is the generating function for the component size distribution $\{P_n\}$, and normalization of probability requires $g(1) = 1$. Indeed, this corresponds to a fixed point of the recursions for $z = 1$ in all models. However, it is a physical solution only if it corresponds to a stable fixed point of the associated recursive equations. Where it fails to be stable, a competing solution with $g(1) < 1$ will take over, yielding a *probability deficit* of magnitude $1 - g(1)$.

For the asymptotically resulting branching process this can be interpreted as being due to a finite probability to obtain an infinite tree. For a finite but large graph, it corresponds to the appearance of a *giant component*, asymptotically containing a finite fraction $1 - g(1)$ of the nodes, and the transition where the naive fixed point loses stability defines a *percolation threshold*, typically of second order. Below the threshold, all components are small, and above it there is a single giant, while the remaining components are small.

4.4.1. The giant in RG

For $z = 1$, the recursion (9) for $g = g(1)$ simplifies to $g \rightarrow e^{c(g-1)}$, with a solution satisfying $cg e^{-cg} = ce^{-c}$. The stability of a solution depends on the magnitude of the Jacobian, given by $ce^{c(g-1)}$, which equals cg when g is a solution.

It has the trivial solution $cg = c$, i.e. $g = 1$. For c smaller than a critical value, $c = 1$, this solution is indeed stable under iteration of the recursion, with the Jacobian given by c .

For $c > 1$, it fails to be stable, and so we must look for another fixed point as the physical solution. Indeed, such a fixed point exists, as follows from looking at a graph of the function $c \rightarrow ce^{-c}$, which has a unique maximum for $c = 1$. Thus, for each $c > 1$ there is a dual value $\hat{c} < 1$ with the same value of this function, yielding the stable solution $g(1) = \hat{c}/c < 1$.

Thus, we have established a probability deficit for $c > 1$, reflecting the existence of a giant component, asymptotically containing a finite fraction $1 - \hat{c}/c$ of the nodes. The critical point $c = 1 \Rightarrow \hat{c} = 1$ defines the percolation threshold, above which there is a finite probability for an arbitrary pair of nodes to be connected via a finite path.

4.4.2. The giant in IRG

For an IRG model, $g = g(1)$ is given by the linear combination $g = \mathbf{r} \cdot \mathbf{g} \equiv \sum_a r_a g_a$, with $\mathbf{g} = \mathbf{g}(1)$ satisfying the recursion $g_a \rightarrow \exp[\sum_b c_{ab} r_b (g_b - 1)]$, as follows from setting $z = 1$ in eqs. (10,11). The stability of the trivial solution $\mathbf{g} = \mathbf{1}$ depends on the spectrum of the local Jacobian matrix $\mathbf{J} = \{c_{ab} r_b\}$.

The case of the largest eigenvalue of \mathbf{J} being exactly unity defines a critical hypersurface in parameter space, beyond which the trivial fixed point $\mathbf{g}(1) = \mathbf{1} \Rightarrow g(1) = 1$ loses stability, and a competing fixed point appears with $g_a(1) < 1 \Rightarrow g(1) < 1$. Again, the corresponding probability deficit $1 - g(1)$ is taken as the probability for winding up in a giant component of size $N(1 - g(1))$.

4.4.3. The giant in DRG

For a DRG model, setting $z = 1$ in eqs. (12,13), yields for $g = g(1)$ and $h = h(1)$ the relation $g = H(h)$ and the recursion $h \rightarrow H'(h)/H'(1)$, with the trivial fixed point $h = 1 \Rightarrow g = 1$. The stability of this is governed by the Jacobian $H''(1)/H'(1) = \langle m(m-1) \rangle / \langle m \rangle$. Stability results if this is smaller than unity, i.e. if $\langle m(m-2) \rangle < 0$, defining the subcritical domain of DRG [4].

In the supercritical domain, there will be a unique competing solution $h < 1$, satisfying $hH'(1) = H'(h)$, yielding $g < 1$, with the corresponding probability deficit indicating the existence of a giant component.

4.4.4. The giant in CDRG

Similarly, in a CDRG model, we can pinpoint the subcritical region by analyzing the stability of the trivial solution $\mathbf{h}(1) = \mathbf{1}$ of the recursion (15) with $z = 1$, amounting to $\mathbf{h} \rightarrow \mathbf{T}\partial H(\mathbf{h})$. The Jacobian amounts to $\mathbf{J} = \mathbf{T}\mathbf{E}$, i.e. the matrix product of \mathbf{T} and the matrix $\mathbf{E} = \{E_{ab}\}$ of second order multivariate combinatorial moments of the colored degree distribution, as defined in eq. (7), and subcriticality corresponds to the largest eigenvalue of \mathbf{J} being smaller than unity.

In the supercritical region, we will have non-trivial solution yielding $g(1) < 1$, with an associated probability deficit and a giant component of corresponding relative size.

5. Discussion

All of the models discussed in this article admit versions with or without the restriction to simple graphs. They share the existence of several nice properties, such as the computability of interesting local and global characteristics, and the existence of a phase transition in the form of a percolation threshold, where a giant component appears.

The sparse RG model is a mathematically very interesting object. Nevertheless, it is severely limited as a model of real-world networks. Its degree distribution is restricted to be Poissonian, and it suffers from a fundamental lack of correlations between edges. Its main importance is as a role model for more general random graph models.

The DRG approach yields a general class of random graph models, and contains a non-simple version of RG as a special case. Although it admits arbitrary degree distributions, it shares with RG a fundamental lack of edge correlations.

Generalizing RG by adding hidden variables in the form of unobservable vertex colors, allowed to affect edge probabilities, yields another general class of models – IRG. It admits arbitrarily many distinct models for a single degree distribution, and displays non-trivial edge correlation. Its most serious limitation lies in the restriction of the degree distribution to a Poissonian mix (which however does not exclude power-behaviour!). It trivially contains RG as a special case, and its restriction to a rank one preference matrix, $c_{ab} = C_a C_b$, defines a class of uncorrelated models, that has been shown to be asymptotically equivalent to the restriction of DRG to Poissonian mixtures [6]. Thus, IRG and DRG define distinct superclasses of the classic RG model, and one might expect that there exists a larger class that contains them both as distinct restrictions.

Such a unified class of models indeed exist. The generalization of DRG to models with unobservable color on individual stubs, that is allowed to

affect the edge probabilities as emerging from the stub pairing statistics, yields a very general class of models – CDRG. It allows for arbitrary degree distributions, as well as for non-trivial edge correlations. It contains as distinct subclasses both DRG (trivially) and IRG (as the restriction of the colored degree distribution to a mix of multivariate Poissonians) [7, 8].

CDRG shares with DRG an interesting relation to Feynman graphs of simple field theories; work is in progress to explore this relation. CDRG should also admit a straightforward extension to cover models also of directed graphs.

A unifying formalism for random graphs appears to be a prerequisite for the possibility to devise a systematic model inference scheme based on the observed properties of real-world networks. CDRG appears to be a step on the way to such a formalism for sparse, truly random graphs.

Acknowledgment

The author thanks the organizers for bringing together a highly interesting workshop, and for the extra bonus of being able to experience the beautiful environment offered by central Krakow.

This work was in part supported by the Swedish Foundation for Strategic Research.

REFERENCES

- [1] P. Erdős, A. Rényi, *Publ. Math. Inst. Hungar. Acad. Sci.* **5**, 17 (1960).
- [2] E. A. Bender, E. R. Canfield, *J. Combinat. Theory, Ser. A* **24**, 296–307 (1978).
- [3] T. Łuczak, in *Poznań, 1989* (A.M. Frieze, T. Łuczak, eds.), *Random Graphs*, vol. 2 (John Wiley & Sons, New York, 1992), pp. 165–182.
- [4] M. Molloy, B. Reed, *Combinat. Prob. Comput.* **7**, 295–306 (1998).
- [5] M. E. J. Newman, S. H. Strogatz, D. J. Watts, *Phys. Rev. E* **64**, 026118 (2001).
- [6] B. Söderberg, *Phys. Rev. E* **66**, 066121 (2002).
- [7] B. Söderberg, *Phys. Rev. E* **68**, 015102(R) (2003).
- [8] B. Söderberg, *Phys. Rev. E* **68**, 026107 (2003).
- [9] B. Bollobás, *Random graphs, 2nd ed.*, Cambridge University Press, Cambridge, 2001.
- [10] S. Janson, T. Łuczak, A. Ruciński, *Random graphs*, Wiley & Sons, New York, 2000.
- [11] P. Flajolet, D. E. Knuth, B. Pittel, *Discr. Math.* **75**, 167–215 (1989).
- [12] Z. Burda, J. D. Correia, A. Krzywicki, *Phys. Rev. E* **64**, 046118 (2001).

- [13] S. N. Dorogovtsev, J. F. F. Mendes, A. N. Samukhin, cond-mat/0204111, 2002.
- [14] M. E. J. Newman, *Phys. Rev. E* **67**, 026126 (2003).