

Component evolution in random intersection graphs

Michael Behrisch *

Institut für Informatik, Humboldt-Universität zu Berlin, 10099 Berlin, Germany

January 21, 2005

We study the evolution of the size of the largest and the second largest component in the random intersection graph model which is suited to reflect the transitivity (or clustering property) visible in real-world networks. We show that certain random intersection graphs differ from $G_{n,p}$ in that they have only a polynomial jump in the evolution of the size of the largest component. On the other hand the moment for the jump is still at the point where the expected vertex degree becomes 1 which is similar to $G_{n,p}$. We also describe a test of our result on a protein network.

1 Introduction

The classical random graph model (introduced by Erdős and Rényi in the early 1960s) considers a fixed set of n vertices and edges that exist with a certain probability $p = p(n)$, independently from each other. It was shown to be inappropriate for describing real-world networks because it lacks certain features of those (e.g. scale free degree distribution and clustering). One of the underlying reasons that are responsible for this mismatch is precisely the independence of the edges, in other words the missing transitivity. In a real-world network, relations between vertices x and y on the one hand and between vertices y and z on the other hand suggest a connection of some sort between vertices x and z .

An *intersection graph* is a graph on vertex set V where each vertex has a subset of a ground set W assigned and two vertices are adjacent if and only if the assigned sets have a non-empty intersection.

We call the ground set W from which the assigned sets are chosen *universal feature set* and its elements *features*. Furthermore the set of vertices V_w holding a specified feature w (which obviously forms a clique) is called *feature clique* while W_v shall denote the set of features assigned to vertex v .

*supported by the DFG research center "Mathematics for key technologies" (FZT 86) in Berlin.

Examples for intersection graphs are the well studied interval graphs on the real line, in this paper however we will only consider finite sets.

A random intersection graph on n vertices with a universal feature set of size m is one where each vertex chooses each feature independently with probability p . A sample of this probability space is denoted by $G_{n,m,p}$.

We consider now and in the following at $m := n^\alpha$ with either $\alpha > 1$ or $0 < \alpha < 1$.

It is sometimes convenient to look at the random intersection graph as a random bipartite graph with bipartition (V, W) and edges occurring between the two classes independently with probability p . A sample from this space will be noted by $B_{n,m,p}$. Note that each component in this graph (ignoring isolated vertices in W) corresponds to a component in the intersection graph and vice versa.

This random model was invented and studied with respect to subgraph appearance by Karoński, Scheinerman and Singer-Cohen in [10], with respect to equivalence to $G_{n,p}$ by Fill, Scheinerman, Singer-Cohen in [6] and with respect to vertex degree distribution by Stark [13]. An algorithmic reconstruction of the feature structure with only the intersection graph as input was given by Behrisch and Taraz in [2]. The first two results and some results concerning connectivity and cliques can also be found in Singer [12].

Extensions to the model were proposed by Godehardt and Jaworski in [7], who modify the distribution of the sizes of the feature cliques and practical relevance of the model was studied by Newman, Strogatz and Watts in [11] and by Guillaume and Latapy in [8].

The aim of this paper is to study the evolution of the largest component in this model. Since components are natural candidates for clusters in graphs it is straightforward to analyze their growth in our random model, thereby getting insight into structural peculiarities of the real-world networks. The component structure for $G_{n,p}$ was already studied by Erdős and Rényi in [5] and there are also results for some models for real-world networks by Chung and Lu [4] and Bollobás and Riordan [3].

The paper is organized as follows. In the next section we describe our results and compare it with the growth of the giant component in $G_{n,p}$. Section 3 states some results on branching processes which will be used for the proofs of the results in Section 5 and 4. We close with some studies on the evolution of a real-world network.

2 The results

Let denote $L_1(G)$ the size of the largest component of G and $L_2(G)$ the size of the second largest component. Our main theorems are:

Theorem 1. *Let $G_{n,m,p}$ be a random intersection graph with $m := n^\alpha$, $\alpha > 1$ and $p^2 m = \frac{c}{n}$. Then for $c < 1$ a.a.s. $L_1(G_{n,m,p}) \leq \frac{9}{(1-c)^2} \ln n$ vertices.*

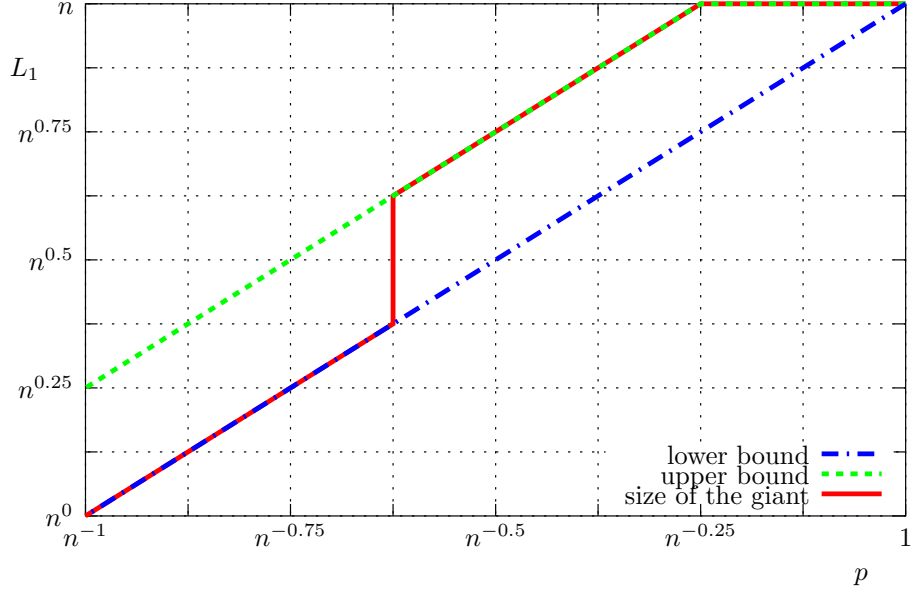


Figure 1: Evolution of the largest component for $\alpha = 0.25$.

Theorem 2. Let $G_{n,m,p}$ be a random intersection graph with $m := n^\alpha$, $\alpha > 1$ and $p^2m = \frac{c}{n}$. Then for $c > 1$ a.a.s. $L_1(G_{n,m,p}) = (1 + o(1))\gamma n$ for an appropriate constant $\gamma = \gamma(c)$ and $L_2(G_{n,m,p}) \leq \frac{50c}{(c-1)^2} \ln n$.

Theorem 3. Let $G_{n,m,p}$ be a random intersection graph with $m := n^\alpha$ and $\alpha < 1$ and $p^2m = \frac{c}{n}$. Then for $c < 1$ a.a.s. $L_1(G_{n,m,p}) \leq \frac{10\sqrt{c}}{(1-c)^2} \sqrt{\frac{n}{m}} \ln m$.

Theorem 4. Let $G_{n,m,p}$ be a random intersection graph with $m := n^\alpha$ and $\alpha < 1$ and $p^2m = \frac{c}{n}$. Then for $c > 1$ a.a.s. $L_1(G_{n,m,p}) = (1 + o(1))\delta\sqrt{mn}$ for an appropriate constant $\delta = \delta(c)$.

As already proven in [12] the "edge probability" p' (meaning the ratio between present edges and all possible edges) in the random intersection graph is closely concentrated around p^2m . Thus the two results above show that for $\alpha > 1$ the largest component in the intersection graph exhibits a jump from logarithmic size to linear size at $p' = \frac{1}{n}$ which is similar to the $G_{n,p}$ behavior. This is also the moment at which in both models the expected degree of a vertex gets larger than 1.

For $\alpha < 1$ the jump is still at the same position but L_1 differs only polynomially as is shown in Figure 1 for $\alpha = 0.25$.

This figure also shows that the size of the largest component jumps approximately from the size of a single feature clique (which is concentrated around pn , see Lemma 9) as a trivial lower bound to the size of the largest component to approximately the sum of the sizes of all feature cliques (which is for the same reasons concentrated around pmn) which is an upper bound to L_1 .

Furthermore we can give a more precise estimation of the size of the largest component for $\alpha < \frac{1}{2}$ at the early stages of the evolution.

Theorem 5. *Let $G_{n,m,p}$ be a random intersection graph with $m := n^\alpha$ and $\alpha < \frac{1}{2}$ and $\ln n \ll pn \ll \frac{\sqrt{n}}{m}$. Then all components are a.a.s. either (feature) cliques or isolated vertices and thus a.a.s. $L_1(G_{n,m,p}) = L_2(G_{n,m,p}) = (1 + o(1))pn$.*

3 Branching processes

In order to discover components in a graph we will use branching processes (for an overview of the topic of branching processes and for references to proofs see [1]) similar to the proofs in Chapter 5 of [9]. We will explore the component by starting at a single vertex, generating its neighbors as descendants in a branching process and then the second neighborhood as their descendants and so forth. As long as the component is still small (which will be made precise later) there is no need to distinguish whether a vertex was already visited in an earlier step of the exploration or not.

Let the random variable X denote the number of neighbors of an arbitrary vertex. The Galton-Watson branching process on the variable X has the following properties:

Theorem 6 (Theorem 5.1 in [9]). *The branching process dies out with probability 1 if $\mathbb{E}[X] \leq 1$, unless $\mathbb{P}[X = 1] = 1$. For $\mathbb{E}[X] > 1$ and $\mathbb{P}[X = 0] > 0$ the probability of extinction is the unique solution of*

$$\sum_{i \geq 0} x^i \mathbb{P}[X = i] = x$$

in the interval $(0, 1)$.

Corollary 7 (Example 5.2 and 5.3 in [9]). *Let X be a random variable with binomial distribution $\text{Bi}(n, p)$ and $np \xrightarrow{n \rightarrow \infty} c > 1$. Then the branching process dies out with probability $\rho(n, p) = 1 - \gamma(c)$ where $\gamma(c)$ is the unique solution of*

$$\gamma + e^{-\gamma c} = 1$$

in the interval $(0, 1)$.

Thus the main thing to do is to estimate the expected value and to overcome the limitations of the branching process which deals with an essentially unbounded domain in contrast to the limited number of vertices in the graph.

The discovery of neighbors is (in contrast to the process used in the $G_{n,p}$ model) a two step process. First we let the vertex in question discover its features and then we let the features find the vertices which hold them. The features used in each step will be ignored in the further process which will

slightly downsize the universal feature set. As we will see later this deviation will not affect the ongoing process very much.

The following estimates are used without proof:

$$(1 - a)^b = (1 + o(1))(1 - ab) \quad \text{for } 0 < a < 1, ab \rightarrow 0 \quad (1)$$

$$e^{-2a} \leq 1 - a \leq e^{-a} \quad \text{for } 0 \leq a \leq \frac{1}{2} \quad (2)$$

Let X be a non-negative random variable with expectation $\mu = \mathbb{E}[X]$ and variance $\text{Var}[X]$. As a special case of Markov's inequality the first moment method states that

$$\mathbb{P}[X \geq 1] \leq \mu. \quad (3)$$

and the second moment method (special case of Tschebyscheff's inequality) that

$$\mathbb{P}[X = 0] \leq \text{Var}[X]/\mu^2 = \frac{\mathbb{E}[X^2]}{\mu^2} - 1. \quad (4)$$

If X is binomially distributed random variable (n trials, each with probability p), then $\mu = np$ and we shall use the following variants of Chernoff's inequality (see Section 2 in [9]):

$$\mathbb{P}[X \geq \mu + t] \leq \exp\left(-\frac{t^2}{2(\mu + t/3)}\right) \quad \text{for } t \geq 0, \quad (5)$$

$$\mathbb{P}[X \leq \mu - t] \leq \exp\left(-\frac{t^2}{2\mu}\right) \quad \text{for } t \geq 0, \quad (6)$$

4 The evolution for $\alpha > 1$

4.1 The size of the feature-set

In evaluating the expected size of the neighborhood of a vertex the following lemma on the size of the feature set of a vertex is a valuable helper.

Lemma 8. *Let v be a fixed vertex in a random intersection graph $G_{n,m,p}$ with $pn = o(1)$ and $p^2mn = \Theta(1)$. Furthermore let $W' \subseteq W$ be a subset of the universal feature set of size at least $m - 2pmn$ and $X_v := |W_v \cap W'|$ denote the random variable counting the number of features of v in W' . Then X_v is very likely close to its expectation or precisely:*

$$\mathbb{P}\left[|X_v - pm| > (pm)^{\frac{3}{4}}\right] \leq \exp\left(-\frac{(pm)^{\frac{1}{2}}}{3}\right)$$

Proof. For the expected number of features selected in W' we have $\mu := \mathbb{E}[X_v] \geq p(m - 2pmn) = pm - O(1)$ and $\mu \leq pm$.

Since the features are selected independently uniformly at random we can use Chernoff inequalities (5) and (6) to bound the deviation from the

expected size.

$$\begin{aligned}
\mathbb{P} \left[Y \geq pm + (pm)^{\frac{3}{4}} \right] &\leq \mathbb{P} \left[Y \geq \mu + (pm)^{\frac{3}{4}} \right] \\
&\leq \exp \left(-\frac{(pm)^{\frac{3}{2}}}{2 \left(\mu + (pm)^{\frac{3}{4}}/3 \right)} \right) \\
&\leq \exp \left(-\frac{(pm)^{\frac{3}{2}}}{2 \left(pm + (pm)^{\frac{3}{4}}/3 \right)} \right) \\
&\leq \frac{1}{2} \exp \left(-\frac{(pm)^{\frac{1}{2}}}{3} \right)
\end{aligned}$$

And for the lower tail using (6):

$$\begin{aligned}
\mathbb{P} \left[Y \leq pm - (pm)^{\frac{3}{4}} \right] &= \mathbb{P} \left[Y \geq \mu + O(1) - (pm)^{\frac{3}{4}} \right] \\
&\leq \exp \left(-\frac{\left((pm)^{\frac{3}{4}} - O(1) \right)^2}{2(pm - O(1))} \right) \\
&\leq \frac{1}{2} \exp \left(-\frac{(pm)^{\frac{1}{2}}}{3} \right)
\end{aligned}$$

Notice that these calculations (and thus the probability for the tails) remain valid even if we remove no features at all.

From the two tails above we may easily conclude the statement of the lemma. \square

4.2 Proofs of Theorem 1 and Theorem 2

Now the proof of Theorem 1 is in reach

Proof of Theorem 1. We prove that for $c < 1$ the branching process starting at an arbitrary vertex v discovering all the vertices one by one will finish in at most $\frac{9 \ln n}{(1-c)^2}$ steps.

From Lemma 8 we know that there is with high probability no large deviation from the expected value in the size of a feature set. Our branching process starting at v now proceeds as follows. At first v discovers its features. If there are too many or too few of them (in the sense of Lemma 8) we abort.

Otherwise we let the features discover the vertices which hold them. Since the feature set of v has size $(1 + o(1))pm$ the probability for an individual vertex w to hold at least one feature in this set is

$$\mathbb{P} [\{v, w\} \in E(G_{n,m,p})] = 1 - (1-p)^{(1+o(1))pm} \stackrel{(1)}{=} (1+o(1))p^2m$$

and the neighbors of v will be chosen independently with this probability. Thus the expected number of new neighbors discovered will be:

$$\mathbb{E}[d(v)] \leq n(1 + o(1))p^2m$$

Now we remove W_v (the feature set of v) from the universal feature set and continue with discovering the features of the neighbors of v the same way we discovered the features of v and so on. We do this at most n times (only n vertices available) thus the probability that we will abort at any step because of the wrong size of the feature set is (due to Lemma 8) bounded by

$$n \exp\left(-\frac{(pm)^{\frac{1}{2}}}{3}\right) \xrightarrow{n \rightarrow \infty} 0.$$

Furthermore we did remove at most $n(1 + o(1))pm < 2pmn$ features from the universal feature set thus Lemma 8 was applicable all the time.

Observe that the probability that v is in a component of size at least k is bounded by the probability that the sum of the degrees of k vertices discovered in the process is at least $k - 1$. Since all features were discovered independent from earlier ones and thus all vertices were discovered in an independent manner, the probability for a component of size at least $k \geq \frac{9 \ln n}{(1-c)^2}$ can be bounded using a Chernoff inequality again. Let Y_i denote the number of neighbors of the i th vertex discovered in the process and notice that the expected value for the sum over the Y_i is bound from above by $(1 + o(1))kp^2mn \leq kc'$ for $c' := \frac{c+1}{2}$.

$$\begin{aligned} n\mathbb{P}\left[\sum_{i=1}^k Y_i \geq k - 1\right] &= n\mathbb{P}\left[\sum_{i=1}^k Y_i \geq kc' + (1 - c')k - 1\right] \\ &\leq n \exp\left(-\frac{((1 - c')k - 1)^2}{2(c'k + (1 - c')k/3)}\right) \\ &\leq n \exp\left(-\frac{(1 - c')^2}{2}k\right). \end{aligned}$$

Resubstituting c' and k shows that this term tends to 0 as n tends to infinity which proves by (3) the theorem. \square

For the appearance of a giant component when $c > 1$ the proof is not as short but it follows again the route of Janson, Łuczak and Rućinski [9].

Proof of Theorem 2. For $c > 1$ we start by proving that there is a.a.s. no component which has less than $k_- := \frac{50c}{(c-1)^2} \ln n$ or more than $k_+ := n^{2/3}$ vertices. This will be the result of examination of the same branching process as in the last proof. We will prove the even harder result that the process whenever it has passed the k_- mark has a.a.s. $\frac{(c-1)}{2}k$ vertices which are to be examined (got discovered as neighbors but were not examined themselves). Notice that in order to prove this we have to look at no more than $k + \frac{c-1}{2}k = \frac{c+1}{2}k$ vertices.

Because of this we exclude in each step at most $\frac{c+1}{2}k_+$ vertices from the further process. Furthermore we do still downsize the universal feature set only for a very small amount for each vertex which discovers its neighbors as in the proof of Theorem 1. This gives independence for all steps of the branching process and thus one can bound the number of neighbors a vertex discovers from below by independent random variables $Y_i^* \in \text{Bi}(n - \frac{c+1}{2}k_+, p'^2m)$ with p' such that $p'^2mn = \frac{3c+1}{4}$. The value for p' results from the lower bound on the size of feature set given by Lemma 8.

Now we can bound the probability of dying out after k steps or having too few discovered (but unexamined) vertices by the probability that

$$\sum_{i=1}^k Y_i^* \leq k - 1 + \frac{c-1}{2}k$$

Now the existence of such a process can be bound by Chernoff inequality (6) and we get with $\mu := \mathbb{E} \left[\sum_{i=1}^k Y_i^* \right] = \frac{3c+1}{4}k - o(k)$ for $k_- \leq k \leq k_+$ and n large enough:

$$\begin{aligned} n \sum_{k=k_-}^{k_+} \mathbb{P} \left[\sum_{i=1}^k Y_i^* \leq k - 1 + \frac{c-1}{2}k \right] &= n \sum_{k=k_-}^{k_+} \mathbb{P} \left[\sum_{i=1}^k Y_i^* \leq \mu - \left(\frac{c-1}{4}k - o(k) + 1 \right) \right] \\ &\leq n \sum_{k=k_-}^{k_+} \exp \left(\frac{- \left(\frac{c-1}{4}k - o(k) + 1 \right)^2}{\frac{3c+1}{2}k} \right) \\ &\leq n \sum_{k=k_-}^{k_+} \exp \left(\frac{- \left(\frac{c-1}{4} \right)^2 k}{3c} \right) \\ &\leq nk_+ \exp \left(\frac{- \left(\frac{c-1}{4} \right)^2 k_-}{3c} \right) \end{aligned}$$

Because of the values for k_- and k_+ given at the beginning of the proof this tends to 0 as n tends to infinity and thus by (3) there is a.a.s. no process stopping between k_- and k_+ .

Now if we have two vertices belonging to components V_1 and V_2 of size at least k_+ the probability that their components are different is equal to the probability that all the vertices in V_1 have not chosen any of the features the vertices in V_2 have selected. According to Lemma 8 the number of (different) features chosen by vertices in V_1 is a.a.s. bounded from below by $k_+ \frac{pm}{2}$. The probability that none of these features was chosen by any of the vertices in V_2 is:

$$(1-p)^{k_+^2 \frac{pm}{2}} \stackrel{(2)}{\leq} \exp \left(-k_+^2 \frac{p^2m}{2} \right) = \exp \left(-n^{\frac{4}{3}} \frac{c}{2n} \right) \xrightarrow{n \rightarrow \infty} 0$$

Now we have that there is a.a.s. only one component with at least k_+ vertices and it remains to show that this component has in fact linear size. In

order to do so we estimate expectation and variance of the random variable Y , denoting the number of vertices in the "small" components of size at most k_- . Let for each vertex $i \in V$ Y_i be the indicator variable for being in a small component.

For a single vertex the probability of being in a small component can be bound from above by the extinction probability of a branching process with distribution $\text{Bi}(n - k_-, (1 - o(1))p^2m)$ and from below by the extinction probability of a branching process with distribution $\text{Bi}(n, (1 + o(1))p^2m)$. The $o(1)$ terms in the two cases bound the possible deviations in the size of feature sets according to Lemma 8.

By Corollary 7 we know that the probability of extinction ρ of these two processes is given by $1 - \gamma(c)$ which results by linearity of expectation into $\mathbb{E}[Y] = (1 - \gamma(c))n$.

In order to have this size a.a.s., we calculate the variance, or precisely using (4) we show that $\mathbb{E}[Y^2] = (1 + o(1))\mathbb{E}[Y]^2$. Two vertices being simultaneously in a small component is an event which occurs either if they are in the same component in that case the probability can be bound by the extinction probability for this component or they are in two components which means two extinctions have to occur independently.

$$\begin{aligned} \mathbb{E}[Y^2] &= \mathbb{E}\left[\left(\sum_{i=1}^n Y_i\right)^2\right] = \sum_{i,j} \mathbb{E}[Y_i Y_j] \\ &\leq n\rho(n,p)k_- + n\rho(n,p)n\rho(n - k_-, p) \\ &= (1 + o(1))n^2\rho(n,p)^2 = (1 + o(1))\mathbb{E}[Y]^2 \end{aligned}$$

By Tschebyscheff's inequality we can conclude that the number of small vertices is a.a.s. $(1 - \gamma(c))n$ hence the largest component is of size $\gamma(c)n$. \square

5 The evolution for $\alpha < 1$

Let $G_{n,m,p}$ be a random intersection graph. Then the probability that there is a feature clique which deviates much from its expected size is exponentially small, or precisely

Lemma 9. *Let $X_w := |V_w|$ be the random variable counting the number of vertices of a fixed feature w in a random intersection graph $G_{n,m,p}$ with $m := n^\alpha$ and $\alpha < 1$. Then:*

$$\mathbb{P}\left[\exists w \in W : |X_w - pn| > (pn)^{\frac{3}{4}}\right] \leq m \exp\left(-\frac{(pn)^{\frac{1}{2}}}{3}\right)$$

Proof. The number of vertices chosen by a feature is a binomially distributed variable. Its deviation from the expected value can therefore be bound by Chernoff inequalities (5) and (6).

First let w be fixed:

$$\begin{aligned}\mathbb{P}\left[X_w > pn + (pn)^{\frac{3}{4}}\right] &\leq \exp\left(-\frac{(pn)^{\frac{3}{2}}}{2(pn + (pn)^{\frac{3}{4}}/3)}\right) \leq \frac{1}{2} \exp\left(-\frac{(pn)^{\frac{1}{2}}}{3}\right) \\ \mathbb{P}\left[X_w < pn - (pn)^{\frac{3}{4}}\right] &\leq \exp\left(-\frac{(pn)^{\frac{3}{2}}}{2pn}\right) \leq \frac{1}{2} \exp\left(-\frac{(pn)^{\frac{1}{2}}}{3}\right).\end{aligned}$$

This results by linearity of expectation (summing over all possible w) and Markov's inequality in

$$\mathbb{P}\left[\exists w \in W : |X_w - pn| > (pn)^{\frac{3}{4}}\right] \leq m \exp\left(-\frac{(pn)^{\frac{1}{2}}}{3}\right).$$

□

If we have an small upper bound for the number of vertices two feature cliques have in common we can simply add the clique sizes (provided we know they are connected) in order to estimate the component size. This bound is achieved by the following lemma.

Lemma 10. *Let Y be the random variable counting the number of vertices having more than one feature in a random intersection graph $G_{n,m,p}$ with $m := n^\alpha$ and $\alpha < 1$. Then for $p^2 m^2 n \gg \ln n$:*

$$\mathbb{P}[Y > 2p^2 m^2 n] \xrightarrow{n \rightarrow \infty} 0$$

and for $p^2 m^2 n \xrightarrow{n \rightarrow \infty} 0$:

$$\mathbb{P}[Y > 0] \xrightarrow{n \rightarrow \infty} 0$$

Proof. For a single fixed vertex v the probability of having more than one feature is (when $pm \rightarrow 0$):

$$\mathbb{P}[|W_v| > 1] = 1 - (1-p)^m - (mp(1-p)^{m-1}) \stackrel{(1)}{=} (1+o(1))m^2 p^2.$$

Since all vertices choose their features independently Y is a binomially distributed variable with expectation $nm^2 p^2$ and the second statement of the lemma follows by Markov inequality. For the first statement we can bound the deviation using Chernoff inequality (5).

$$\mathbb{P}[Y > 2p^2 m^2 n] \leq \mathbb{P}[Y > 2\mathbb{E}[Y]] \leq \exp\left(-\frac{3nm^2 p^2}{8}\right) \xrightarrow{n \rightarrow \infty} 0.$$

□

Now we can start proving the component evolution for $\alpha < 1$.

Proof of Theorem 3. The trick we use here is that we interchange the role of the feature set and the vertex set and look at the largest component in the feature set instead of one in the vertex set. As we know from Theorem 1 there will be no component containing more than $\frac{9}{(1-c)^2} \ln m$ features. Together with lemma 9 we can conclude that the size of the largest component is a.a.s. bound by

$$\frac{9}{(1-c)^2} \ln m \cdot (1+o(1))pn \leq \frac{10\sqrt{c}}{(1-c)^2} \sqrt{\frac{n}{m}} \ln m.$$

□

Proof of Theorem 4. We use the same method as in the last proof. With exactly the same argument we already have a.a.s. an upper bound for the size of the largest component of

$$\gamma\left(\frac{1}{c}\right) \cdot (1+o(1))pn \leq (1+o(1))\sqrt{c}\gamma\left(\frac{1}{c}\right) \sqrt{mn}.$$

The lower bound can be achieved because the size of the component can be bound by the sum over the sizes of all cliques minus the number of vertices which occur in more than one clique multiplied with the multiplicity they occur. Or more precise (with W_L denoting the set of features in the giant component in W and V_L denoting the vertices linked to it):

$$\begin{aligned} |V_L| &= \sum_{w \in W_L} |V_w| - \sum_{v \in V_L, |W_v| > 1} (|W_v| - 1) \\ &\geq \gamma\left(\frac{1}{c}\right) m(1+o(1))pn - \sum_{v \in V_L, |W_v| > 1} \max_{v \in V} \{|W_v|\} \end{aligned}$$

The probability for the existence of a vertex with more than $\ln m$ features is bound by $n(pm)^{\ln m}$ which tends to 0 for our choice of p . Furthermore we know from Lemma 10 that there are at most $2p^2m^2n = 2cm$ vertices with more than one feature. Therefore

$$\begin{aligned} |V_L| &\geq \gamma\left(\frac{1}{c}\right) m(1+o(1))pn - 2cm \ln m \\ &= (1+o(1))\gamma\left(\frac{1}{c}\right) \sqrt{cmn} - 2cm \ln m \\ &= (1+o(1))\gamma\left(\frac{1}{c}\right) \sqrt{cmn} \end{aligned}$$

Setting $\delta = \gamma\left(\frac{1}{c}\right) \sqrt{c}$ this gives the statement of the theorem. □

Proof of Theorem 5. The statement follows directly from Lemma 10 and Lemma 9 because if there are no vertices with more than one feature there are only isolated vertices and feature cliques. □

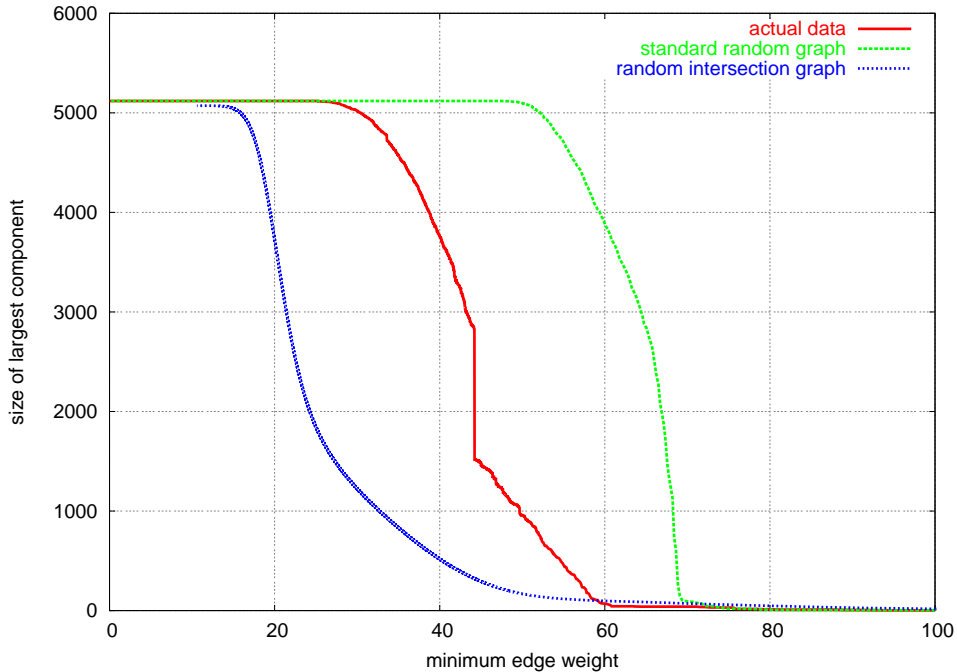


Figure 2: Evolution of the largest component in the protein graph.

6 Experiments

We tested our result on an instance of a complete edge-weighted real world network on 5119 vertices. Here parts of proteins serve as vertices and the edge-weight describes their spatial similarity. If we look at the subgraph of this graph containing all edges with weight greater than a fixed value s (where greater edge weights indicate higher similarity) we can simulate an evolution of this network by gradually decreasing s . Thus first the highly analogue parts get connected and bit by bit also the less similar ones connect to the components.

The evolution found this way differs significantly from a graph in which the same weights are distributed uniformly at random among the edges (see Figure 2).

The most striking difference is the slow growth of the largest component in the stages after it has only very few vertices (minimum edge weight between 40 and 60). A similar behavior cannot be modelled using standard random graphs where L_1 is either logarithmic or linear in the number of vertices. As one can see in Figure 2 the random intersection graph resembles this steady aggregation of vertices to the largest component very well.

References

- [1] K. B. Athreya and A. N. Vidyashankar. Branching processes. Technical report, University of Georgia – Department of Statistics, 1999.
- [2] M. Behrisch and A. Taraz. Efficiently covering complex networks with cliques of similar vertices. Preprint, November 2004.
- [3] B. Bollobás and O. Riordan. Slow emergence of the giant component in the growing m -out graph. to appear in *Random Structures and Algorithms*.
- [4] F. Chung and L. Lu. Connected components in random graphs with given expected degree sequences. *Annals of Combinatorics*, 6:125–145, 2002.
- [5] P. Erdős and A. Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 5:17–61, 1960.
- [6] J. A. Fill, E. R. Scheinerman, and K. B. Singer-Cohen. Random intersection graphs when $m = \omega(n)$: An equivalence theorem relating the evolution of the $G(n, m, p)$ and $G(n, p)$ models. *Random Structures and Algorithms*, 16(2):156–176, March 2000.
- [7] E. Godehardt and J. Jaworski. Two models of random intersection graphs and their applications. *Electronic Notes in Discrete Mathematics*, 10, 2001.
- [8] J.-L. Guillaume and M. Latapy. Bipartite structure of all complex networks. *Information Processing Letters*, 90:215–221, 2004.
- [9] S. Janson, T. Łuczak, and A. Ruciński. *Random Graphs*. John Wiley & Sons, 2000.
- [10] M. Karoński, E. R. Scheinerman, and K. B. Singer-Cohen. On random intersection graphs: The subgraph problem. *Combinatorics, Probability and Computing*, 8:131–159, 1999.
- [11] M. E. J. Newman, S. H. Strogatz, and D. J. Watts. Random graphs with arbitrary degree distributions and their applications. *Physical Review E*, 64, 2001.
- [12] K. B. Singer. *Random Intersection Graphs*. PhD thesis, John Hopkins University, Baltimore, Maryland, 1995.
- [13] D. Stark. The vertex degree distribution of random intersection graphs. *Random Structures and Algorithms*, 24(3):249–258, May 2004.