

# Message passing problem on random graphs

Sebastian M. Krause

## I. QUESTION

Lets assume the generalized configuration model graph ensemble with  $N$  nodes, where each degree sequences  $\{k_i\}$  occurs with probability  $\prod_i p_{k_i}$  with the degree distribution  $p_k$ . (see M.E.J. Newman: Networks, an Introduction; 2010; Eq. (13.30). The generalization of the configuration model is important for numerics with small networks, where many network realizations are sampled for averaging.) Lets additionally assign to every node  $i$  a color  $c_i \in 1, 2, \dots, C$ . The color sequence  $\{c_i\}$  has probability  $\prod_i r_{c_i}$  with the color distribution  $r_c$ . How large is the fraction of node pairs, which can be connected via a set of paths, such that for every color there exists a path avoiding this color?

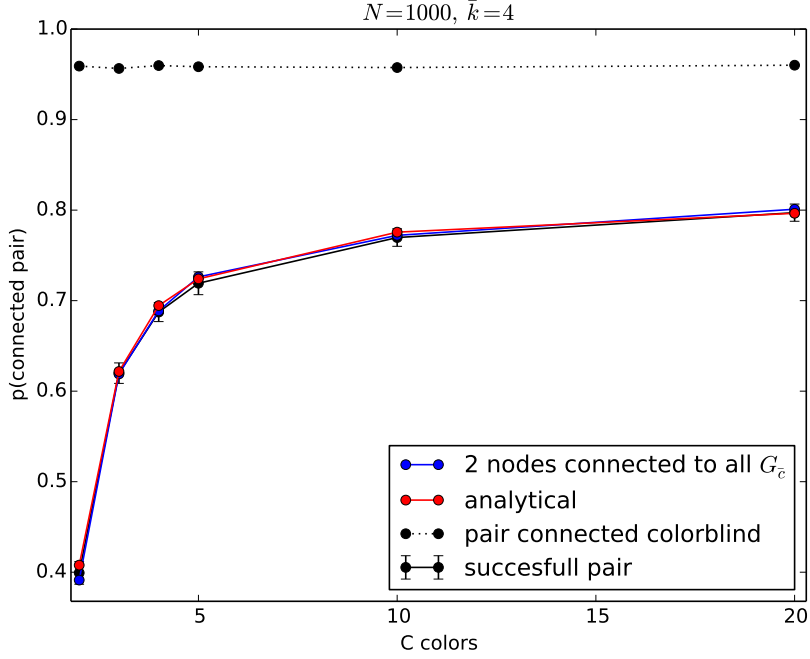


FIG. 1: Fraction of node pairs connected with color avoiding paths as described in the text. Colors are distributed randomly on Poisson graphs with average degree  $\bar{k} = 4$ . Black symbols with error bars show results for networks of size  $N = 1000$ , with samples of 200 node pairs, averaged over 50 networks. Connected pairs are suppressed compared to the existence of colorblind paths, as indicated with black symbols. Blue symbols show the fraction of node pairs where each node is connected to color avoiding giant components, confirming the understanding using percolation theory. Red symbols show analytical results of percolation theory.

For numerical results, we generated Poisson graphs with  $N = 1000$  (we used the Poisson distribution as an approximation to the binomial distribution in ER-graphs), and distributed  $C$  colors over the nodes, with  $r_c = 1/C$ . To check for a single pair, if it can be connected in the way as described above, we removed all nodes with the first color (except for the node pair under consideration) and searched for the shortest path (using a library function returning an empty set, if no path exists). Then we took the original graph and removed all nodes with the second color etc. If for all colors the shortest path exists, we considered the pair as successful. The case of a successful pair with a smaller number of paths, where e.g. one path avoids two colors, is included, as such a path could just be counted twice with our procedure. We repeated this 200 times with node pairs randomly chosen, and calculated the fraction of successful node pairs. See the black straight line in figure 1, with averages over 50 network realizations. The dotted line shows the fraction of node pairs, which are connected with a path which can have all colors.

## II. THEORY

### A. Connection to percolation

For estimating the fraction of successful pairs analytically, results from percolation theory can be used. First of all, the existence of a (colorblind) giant component clearly is a prerequisite for the existence of a macroscopic fraction of successful node pairs. With the generating functions of degree  $g_0(z) = \sum_k p_k z^k$  and excess degree  $g_1(z) = \sum_k q_k z^k$ , the size of the giant component  $S$  can be calculated (assuming infinite networks which are locally treelike) using the average probability  $u$ , “that a vertex is not connected to the giant component via its connection to some particular neighboring vertex” (Newman, page 461):

$$u = g_1(u) \quad (1)$$

$$S = 1 - g_0(u). \quad (2)$$

Lets call the set of all nodes belonging to the giant component as  $\mathcal{G}$ . Another prerequisite is the existence of a giant component after deleting all nodes of one of the colors  $c$ . Lets call the analogue of  $u$  after all nodes of color  $c$  are deleted as  $u_{\bar{c}}$ , the set of all nodes in the remaining giant component as  $\mathcal{G}_{\bar{c}}$ , and its size as  $S_{\bar{c}}$ . We have

$$\phi_{\bar{c}} = 1 - r_c \quad (3)$$

$$u_{\bar{c}} = 1 - \phi_{\bar{c}} + \phi_{\bar{c}} g_1(u_{\bar{c}}) \quad (4)$$

$$S_{\bar{c}} = \phi_{\bar{c}}(1 - g_0(u_{\bar{c}})). \quad (5)$$

Finally a node pair is for sure successful, if for both nodes the following holds: For every color  $c$  there exists at least one neighbor belonging to  $\mathcal{G}_{\bar{c}}$ . Lets call the set of all nodes fulfilling this condition as  $\mathcal{G}_{\text{color}}$ , and its size as  $S_{\text{color}}$ . The fraction of successful node pairs should be approximately  $S_{\text{color}}^2$ .

We tested this hypothesis with numerical results. For every color  $c$ , we labeled all the nodes with a boolean variable, if they belong to  $\mathcal{G}_{\bar{c}}$  (more precise, if they belong to the largest component, as the networks are finite). A single node can belong to different components  $\mathcal{G}_{\bar{c}}$ . Then we searched for nodes which can be successful in a node pair. For a single node, we iterated over all the neighbors and collected all the memberships  $\mathcal{G}_{\bar{c}}$  present. If for all colors  $c$  there is at least one membership  $\mathcal{G}_{\bar{c}}$  among the neighbors, the node belongs to  $\mathcal{G}_{\text{color}}$  and is potentially successful in a pair. We calculated the fraction  $S_{\text{color}}$  over all nodes, and by squaring this we got an estimate for successful node pairs. This procedure ignores paths in small components, but is a good estimate, as can be seen with the blue line in figure 1. This procedure is much faster as well.

### B. Analytical results for the percolation problem

As we have tested numerically that  $S_{\text{color}}$  can be used to describe the success of node pairs in connecting over paths with avoided colors, it is useful to assess this quantity analytically. We will do this assuming an infinite, locally treelike network. We calculate  $S_{\text{color}}$  as the probability, that a randomly chosen node belongs to  $\mathcal{G}_{\text{color}}$ . A randomly chosen node has exactly  $k$  links with probability  $p_k$ . We need the probability that these links connect to all components  $\mathcal{G}_{\bar{c}}$ . As  $\mathcal{G}_{\bar{c}}$  are subsets of  $\mathcal{G}$ , only links connecting to  $\mathcal{G}$  can contribute. We therefore use the conditional probability

$$Q_{k,k'} = \binom{k}{k'} (1-u)^{k'} u^{k-k'} \quad (6)$$

that out of  $k$  links  $k'$  connect to the giant component. Further the distribution of colors among the nodes these links connect to is crucial.

$$R_{k',\vec{\kappa}} = \frac{k'!}{\kappa_1! \times \dots \times \kappa_C!} (r_1)^{\kappa_1} \times \dots \times (r_C)^{\kappa_C} \delta_{k',\kappa_1+\dots+\kappa_C} \quad (7)$$

denotes the conditional probability that out of those  $k'$  links  $\kappa_1$  connect to nodes of color 1,  $\kappa_2$  links connect to nodes of color 2 etc.

Lets now concentrate on one color  $c$ . We have  $k' - \kappa_c$  links which potentially can connect to the desired component  $\mathcal{G}_{\bar{c}}$ . According to the choices we have made so far, those links connect to the giant component  $\mathcal{G}$  and none of the nodes they are connecting to has color  $c$ . Therefore, a single of those links does not connect to  $\mathcal{G}_{\bar{c}}$  with the conditional probability

$$U_{\bar{c}} = 1 - \frac{1 - u_{\bar{c}}}{(1 - u)\phi_{\bar{c}}}. \quad (8)$$

The last term is the probability, that over a single link a connection to  $\mathcal{G}_{\bar{c}}$  is established, if this link already fulfills the following precondition: It connects to  $\mathcal{G}$  and at the same time to a node without color  $c$ . This precondition has probability  $(1-u)\phi_{\bar{c}}$ , as colors are randomly distributed and therefore are not correlated with the probability  $u$  or  $1-u$ . As the links connecting to  $\mathcal{G}_{\bar{c}}$  are a subset of all links fulfilling the precondition, the conditional probability can be calculated by dividing with the probability of the precondition. Notice that the additional information of the explicit color  $c'$ , instead of only stating that the color is not  $c$ , does not alter the results, as a further restriction of the colors would meet the numerator and denominator identically and therefore would cancel out.

There is at least one link connecting to  $\mathcal{G}_{\bar{c}}$  with probability  $1 - (U_{\bar{c}})^{k' - \kappa_c}$ . The success probabilities for different colors have to be multiplied, as all  $\mathcal{G}_{\bar{c}}$  have to be reached at the same time. Putting everything together we have

$$S_{\text{color}} = \sum_{k=0}^{\infty} p_k \sum_{k'=0}^k Q_{k,k'} \sum_{\kappa_1, \dots, \kappa_C=0}^{k'} R_{k', \vec{\kappa}} \prod_{c=1}^C [1 - (U_{\bar{c}})^{k' - \kappa_c}]. \quad (9)$$

Results for Poisson graphs are shown in figure 1 with the red line, showing  $S_{\text{color}}^2$  as the probability of two nodes to be connected via all Components  $\mathcal{G}_{\bar{c}}$  simultaneously. Instead of evaluating the sums over  $k'$  and  $\vec{\kappa}$  in eq. 9, we sampled 5000 events for every  $k$ . The outcome compares well with numerical results.

### C. Limiting case of small color frequencies

In the limit of high numbers of colors  $C$  together with color frequencies  $r_c \rightarrow 0$ , the single paths have to avoid only a small part of nodes. Therefore we expect  $U_{\bar{c}} \rightarrow 0$ : If a link connects to the colorblind giant component, it will almost never fail to connect to the color avoiding component. We can use this idea to find a limiting case, and to compare to standard percolation. We use  $\sum_{k=0}^{\infty} p_k \sum_{k'=0}^k Q_{k,k'} \sum_{\kappa_1, \dots, \kappa_C=0}^{k'} R_{k', \vec{\kappa}} = 1$  (total probability is one) together with the fact that  $\prod_{c=1}^C [1 - (U_{\bar{c}})^{k' - \kappa_c}] = 0$  whenever  $k' < 2$  (at least for one color then  $k' - \kappa_c = 0$ ). Therefore we have

$$S_{\text{color}} < 1 - \sum_{k=0}^{\infty} p_k [u^k + k(1-u)u^{k-1}] \quad (10)$$

$$= 1 - g_0(u) - (1-u) \left. \frac{dg_0(z)}{dz} \right|_{z=u} \quad (11)$$

$$= S_{\text{color}, \infty}. \quad (12)$$

The upper limit  $S_{\text{color}, \infty}$  for  $S_{\text{color}}$  reflects the fact, that every node has to be connected to the giant component at least over two links. It therefore includes a reduction compared to the standard percolation result  $1 - g_0(u)$ . In the limit of many colors and small probabilities  $r_c$ ,  $S_{\text{color}}$  can come close to  $S_{\text{color}, \infty}$ , as in this case  $U_{\bar{c}}$  comes close to zero and only nodes fail which have less than two links connecting to the giant component. This result is closely connected to  $k$ -core percolation with  $k = 2$ . Note that  $k$ -core percolation shows a continuous phase transition for  $k = 2$ , and only for  $k > 2$  has the well known discontinuous behavior.

## III. RESULTS

### A. Poisson graphs

In figure 2 the dependence of  $S_{\text{color}}$  on the average degree is shown for different numbers of colors  $C$  ( $r_c = 1/C$ ). Comparing to the standard giant component size  $S$  (dashed black line in the figure), the percolation sets in at increasing  $\bar{k}$  with smaller numbers of colors, and the component size grows slower to the saturation value of one. The circles show numerical results with  $N = 1000$  and 10 network realizations, the lines show results of equation 9, both correspond well.

The suppression of the number of connected nodes can be understood as a combination of two effects. The first effect is purely topological and can be understood with  $S_{\text{color}, \infty}$  (shown with dashed red line). It means that only nodes can belong to  $S_{\text{color}}$ , which are connected to the colorblind giant component over at least two links. We can confirm that  $S_{\text{color}} \rightarrow S_{\text{color}, \infty}$  for high numbers of colors  $C$  with the results for  $C = 50$ .  $S_{\text{color}, \infty}$  is remarkably reduced compared to  $S$  for small  $k$ , but has the same critical parameter. For the Poisson graph we have  $S_{\text{color}, \infty} = S - \bar{k}S(1-S)$ , and for small positive  $\bar{k} - 1$  the giant component grows approximately with  $S \approx 2(\bar{k} - 1)/\bar{k}^2$ . Therefore

$$S_{\text{color}, \infty} \approx 2(\bar{k} - 1)^2/\bar{k}, \quad (13)$$

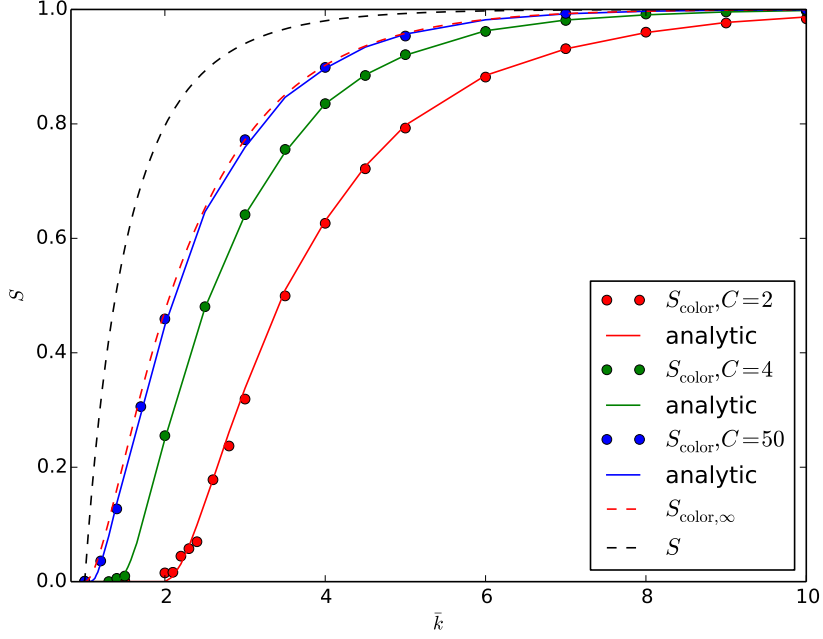


FIG. 2: Dependence of  $S_{\text{color}}$  on the average degree for different number of colors. Symbols show numerical results, the straight lines analytical results. For comparison, the giant component size  $S$  is shown (black dashed).  $S_{\text{color}}$  is reduced due to two mechanisms: First, every node has to be connected to the giant component via two links. The according fraction of nodes  $S_{\text{color},\infty}$  is shown with a red dashed line. Second, increasing color frequencies further decrease  $S_{\text{color}}$ .

which grows slowly for small parameter  $\bar{k} - 1$ .

The second effect is connected to finite color frequencies  $r_c$  which further reduces the percolating fraction of nodes. This also changes the critical value  $\bar{k}_{\text{crit}}$  and the critical exponent  $\beta$ . We will discuss the critical behavior for the more general case of heterogeneous color distributions  $r_c$ .

With general color distributions  $r_c$  ( $\sum_c r_c = 1$ ), the color with the largest probability  $r_c$  dominates the behavior, as it corresponds to the largest conditional link failure probability  $U_{\bar{c}}$  in equation 9. For Poisson graphs,  $U_{\bar{c}}$  falls below one at  $\bar{k} = 1/(1 - r_c)$ , and as long as one  $U_{\bar{c}}$  is one, equation 9 gives zero. Therefor the critical value is

$$\bar{k}_{\text{crit}} = 1/(1 - \max_c r_c). \quad (14)$$

In figure 3 upper left, analytical results for a highest color frequency  $r_1 = 1/3$  are shown. The corresponding critical value is  $\bar{k}_{\text{crit}} = 1/(1 - \max_c r_c) = 3/2$  as expected. Different color distributions were used with different degeneration  $n_{\text{deg}}$  of the highest color frequency and  $C = 10$  colors. We see the same critical value, but different critical exponents  $\beta = n_{\text{deg}}$  for  $S_{\text{color}} \propto (\bar{k} - \bar{k}_{\text{crit}})^\beta$ . This analytical result can be confirmed with numerical results, as shown on the bottom of the figure.

We can understand the critical exponent with expanding equation 9 using  $U_{\bar{c}} = 1 - \varepsilon$  for the highest color frequency. First of all, we have  $(U_{\bar{c}})^{k' - \kappa_c} \approx 1 - (k' - \kappa_c) \times \varepsilon$ , and therefore we find  $\prod_{c=1}^C [1 - (U_{\bar{c}})^{k' - \kappa_c}] \propto \varepsilon^{n_{\text{deg}}}$  if  $k' - \kappa_c > 0$  or  $\prod_{c=1}^C [1 - (U_{\bar{c}})^{k' - \kappa_c}] = 0$  else. As equation 9 is therefore a superposition of either vanishing terms or terms with leading order  $\varepsilon^{n_{\text{deg}}}$ , and we know from standard percolation that  $\varepsilon \propto (\bar{k} - \bar{k}_{\text{crit}})$ , we have

$$\beta = n_{\text{deg}}. \quad (15)$$

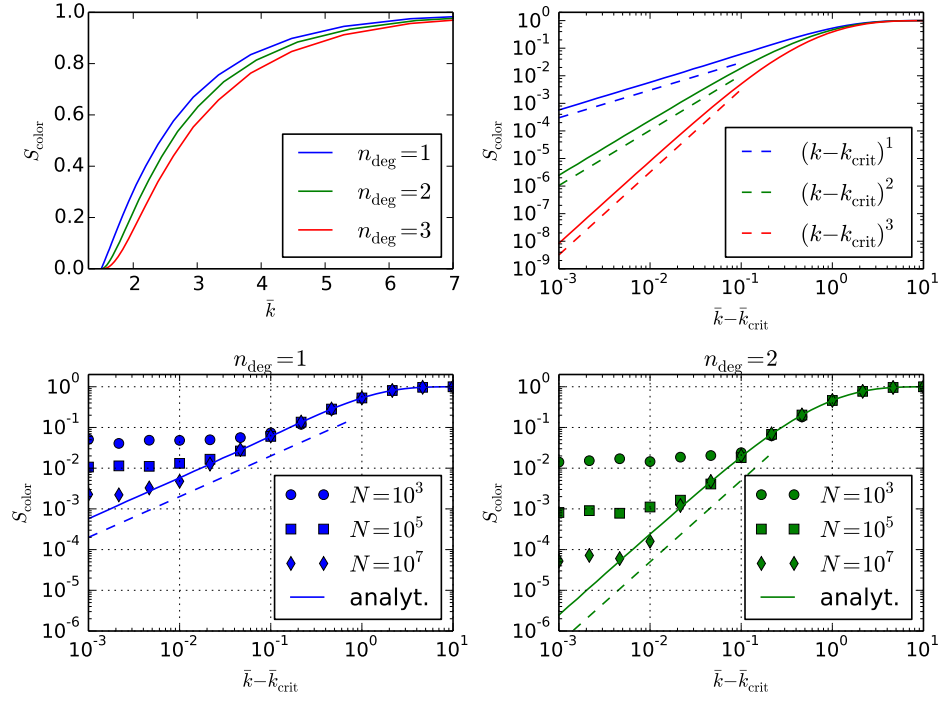


FIG. 3: For heterogeneous color distributions  $r_c$ , the highest frequency determines the behavior. On the upper left, analytical results for Poisson graphs are shown for different color distributions with  $\max_c r_c = 1/3$  and  $C = 10$ , all with the same  $\bar{k}_{\text{crit}} = 1/(1 - \max_c r_c) = 3/2$ . The critical exponent is determined by the degeneration of the highest color frequency. On the upper right, the same results are shown with a log-log plot confirming  $\beta = n_{\text{deg}}$ . This can be understood with an expansion of equation 9. On the bottom, the analytical results are compared to numerical results which converge to the expected critical behavior with system size.

### B. Broad degree distribution

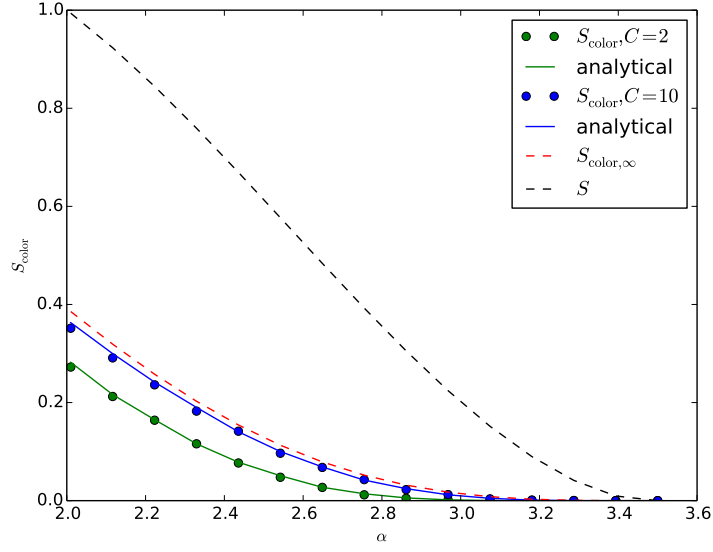


FIG. 4: The same as in figure 2 for scale free degree distributions.

In figure 4, results for graphs with broad degree distributions with  $p_k = nk^{-\alpha}$  are shown.  $n$  is a normalization constant, and the giant component exists between  $\alpha = 2$  and  $\alpha = 3.5$ . We see a strong reduction of  $S_{\text{color},\infty}$  compared to  $S$ , while the number of colors ( $r_c = 1/C$ ) plays a minor role. Numerical results are averages over 50 networks of size  $N = 10\,000$ . For evaluating equation 9, 1000 events were sampled for every  $k$ .

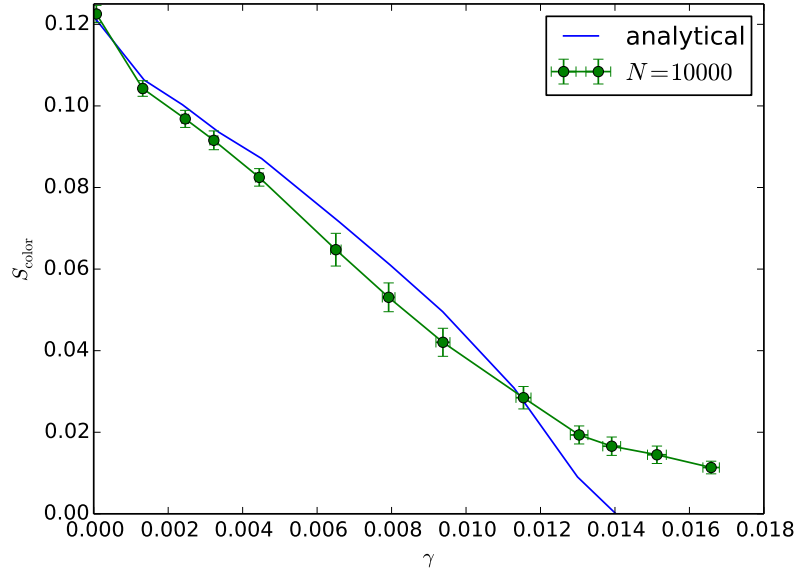


FIG. 5:  $S_{\text{color}}$  drops fast, if a fraction  $\gamma$  of nodes with the highest degree is restricted to one color, while the nodes with smaller degree can have one of two colors. This is shown for a network with  $\alpha = 2.3$  and  $N = 10\,000$  with symbols. Analytical results need a modified version of equation 9 as described in the text. Results are shown with the blue line. Without diversity on the hubs, nodes cannot communicate in the desired way.

For broad degree distributions, color distributions can show an additional type of heterogeneity, as a dependence of frequencies on the degree of a node can strongly influence the behavior. We used two colors, where the first color has a frequency of  $\tilde{r}_{1,k} = 1$  for all degrees  $k \geq k_{\text{step}}$  larger than a certain  $k_{\text{step}}$ . These nodes have a probability of  $\gamma = \sum_{k=k_{\text{step}}}^{\infty} p_k \tilde{r}_{1,k}$ . Accordingly  $\tilde{r}_{2,k} = 0$  for  $k \geq k_{\text{step}}$ , and probabilities for smaller degrees are chosen such that  $\sum_k p_k \tilde{r}_{c,k} = 1/2$ , meaning that every color on average occupies one half of the nodes. Figure 5 shows results for an ensemble with  $\alpha = 2.3$ . Analytical results show that already for a portion of  $\gamma = 1.4\%$  of the largest nodes occupied by the first color exclusively,  $S_{\text{color}}$  vanishes. Numerical results with  $N = 10000$  confirm this behavior, but show finite size effects.

Move the following to appendix and explain in more detail, repeating all equations needed with replacements?. For the analytical results, Equation 9 had to be adapted for the degree dependence of colors: The ingredient  $R_{k',\tilde{r}}$  as given in equation 7 had to be calculated with the replacement  $r_c \rightarrow \sum_k k p_k \tilde{r}_{c,k} / \bar{k}$  which gives the average probability to find a color  $c$  following a random link. The same replacement of  $r_c$  has to be done for  $\phi_{\bar{c}} = 1 - r_c$  in equation 8 for calculating  $U_{\bar{c}}$ . Finally, for calculating  $u_{\bar{c}}$  as needed in the same equation, equation 4 has to be modified:  $u_{\bar{c}} = 1 - f_1(1) + f_1(u_{\bar{c}})$  has to be solved with the modified generating function of excess degree  $f_1(z) = \sum_k q_k (1 - \tilde{r}_{c,k}) z^k$ .

### C. Application: Network of autonomous systems

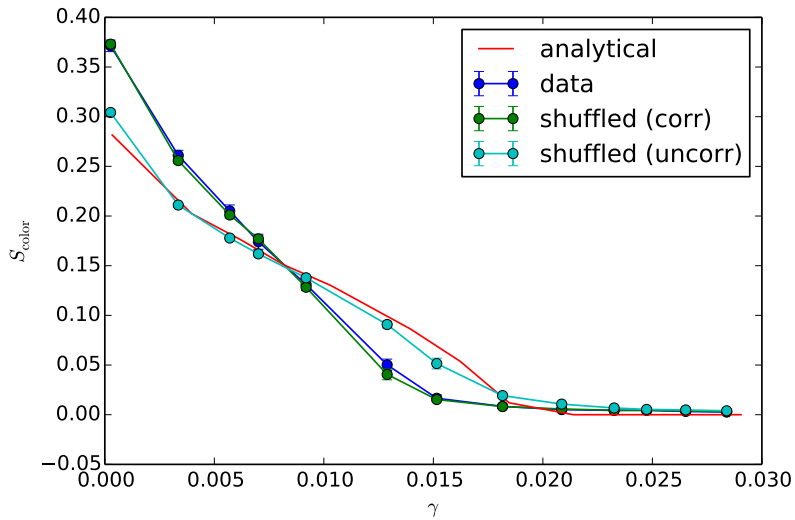


FIG. 6: Blue symbols show results For the network of autonomous systems, where colors are distributed in the same way as in figure 5. Averages were taken over 10 color distributions. Analytical results for an ensemble of infinite size and power law tail ( $\alpha = 2.0304$ ) are shown with the red line and reproduce the results qualitatively. Deviations are due to degree-degree-correlations which are reserved in shuffled networks shown with blue symbols, while results with ignoring correlations are shown with light blue symbols.

The blue symbols in figure 6 show results for the autonomous systems network, where colors were distributed with degree-dependence over the nodes as described at the end of the last section. Averages were taken over 10 realizations of the color distributions. As expected from our results for scale free degree distributions, the  $S_{\text{color}}$  drops to 0 even for small fraction  $\gamma$  which is exclusively of one color. That means that if there is no heterogeneity in the highly connected servers, it is not possible to avoid e.g. software versions. This is also interesting in the following sense: It is known that secret services try to store all decrypted data running through servers they get to decrypt it later. As this is connected to technical afford, the services will more likely monitor the large servers. Therefore it would be beneficial, if once using encryption, to sent parts of the message using small servers. Unfortunately, this seems to be impossible, the services only have to monitor a low percentage of servers to hinder alternative paths.

In order to assess the predictive power of our analytical method, we used a model ensemble. Up to  $k = 10$ , the degree distribution of the autonomous systems was used. For higher degrees, the tail was assumed to be a power law and the exponent was fitted using maximum likelihood to be  $\alpha \approx 2.0304$ . Results for the according ensemble are shown with the red line. The qualitative behavior is represented well. To understand deviations, we compared to data from shuffled networks starting with the original data. Shuffling with ignoring degree-degree-correlations while

only keeping the degree sequence gives results close to the analytical results (light blue symbols). Shuffling with also keeping degree-degree-correlations gives results close to the original network (green symbols). Therefore, deviations between our theory and the data arise mainly due to degree-degree-correlations.

### Appendix: Independence of components for Poisson graphs

In equation 9, the probabilities  $1 - U_{\bar{c}}^{k' - k_c}$  for different colors may show dependencies among each other limiting the usability of the equation. In the following we discuss one particular case for Poisson graphs, in order to illustrate a sort of independence. Lets assume we only have one link connecting to a node in the giant component, and the node has color 3. Can the probability, that this link connects to  $\mathcal{G}_{\bar{1}}$  and  $\mathcal{G}_{\bar{2}}$  at the same time really be written as the product  $(1 - U_{\bar{1}})(1 - U_{\bar{2}})$ ? As for Poisson graphs we have  $S = 1 - u$ , instead of discussing link probabilities  $u$  we can concentrate on node probabilities  $S$ .

We use the notation  $S(\mathcal{Y}|\mathcal{X})$  to denote the fraction of nodes in a set  $\mathcal{Y}$  which is a subset of  $\mathcal{X}$ . With the set  $\mathcal{N}$  of all nodes in the network we have  $S_{\bar{c}} = S(\mathcal{G}_{\bar{c}}|\mathcal{N})$ . Clearly they must be dependent, as  $S(\mathcal{G}_{\bar{1}} \cap \mathcal{G}_{\bar{2}} \cap \dots \cap \mathcal{G}_{\bar{C}}|\mathcal{N}) = 0 \neq S(\mathcal{G}_{\bar{1}}|\mathcal{N}) \times \dots \times S(\mathcal{G}_{\bar{C}}|\mathcal{N})$ . Lets call the set of all nodes with color  $c$  as  $\mathcal{A}_c$ , and the set of all nodes without this color as  $\mathcal{A}_{\bar{c}}$ .

Back to our case with the Poisson graph, we have for  $c = 1, 2$

$$1 - U_{\bar{c}} = \frac{S_{\bar{c}}}{S(1 - r_c)} \quad (16)$$

$$= S(\mathcal{G}_{\bar{c}} \cap \mathcal{A}_{\bar{c}}|\mathcal{G} \cap \mathcal{A}_{\bar{c}}) \quad (17)$$

$$= S(\mathcal{G}_{\bar{c}} \cap \mathcal{A}_3|\mathcal{G} \cap \mathcal{A}_3). \quad (18)$$

The last equation is intuitive, as the coloring of nodes is random, and it was tested numerically (results not shown). In order to test if the product  $(1 - U_{\bar{1}})(1 - U_{\bar{2}})$  reflects the probability of connecting to  $\mathcal{G}_{\bar{1}}$  and  $\mathcal{G}_{\bar{2}}$  at the same time, we have to check if

$$S(\mathcal{G}_{\bar{1}} \cap \mathcal{G}_{\bar{2}} \cap \mathcal{A}_3|\mathcal{G} \cap \mathcal{A}_3) = S(\mathcal{G}_{\bar{1}} \cap \mathcal{A}_3|\mathcal{G} \cap \mathcal{A}_3) \times S(\mathcal{G}_{\bar{2}} \cap \mathcal{A}_3|\mathcal{G} \cap \mathcal{A}_3) \quad (19)$$

holds. The comparison of the left hand side and the right hand side is shown in figure 7 for networks with  $N = 1000$ . For every number of colors  $C = 3, 4, 5, 10$ , ten network realizations were used. So we have tested numerically, that equation 9 is useful at least in this very simple case. This was to illustrate a sort of independence of the conditional probabilities for connecting to color-avoiding giant components. The generalization to many links and general degree distributions is missing so far.



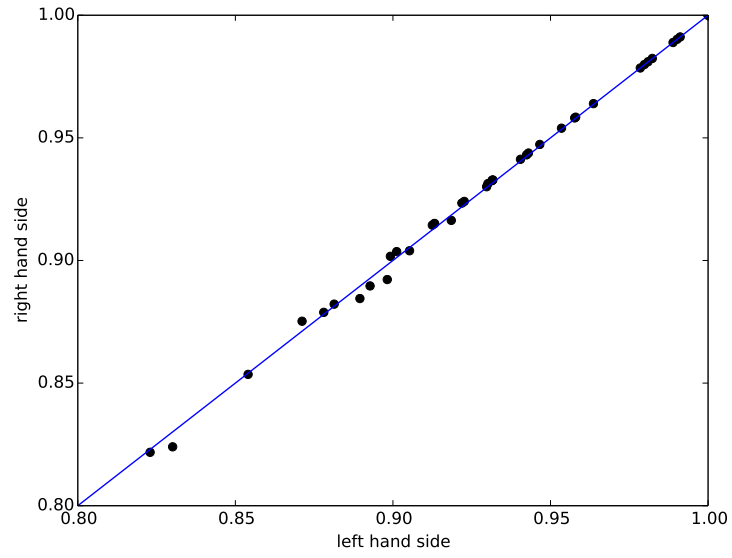


FIG. 7: Calculating both sides of equation 19 for Poisson networks with  $N = 1000$  motivates the usage of the conditional probabilities  $U_{\bar{e}}$  as independent quantities.