

GENERATION OF SYNTHETIC NETWORKS ON DIFFERENT ENSEMBLES AND THE EFFECT OF SELF-LOOPS

This document constitutes the manual for the **Multi Edge Randomizer** software, where details on the implementation of the simulated ensembles are provided. For more details see [1] and [2].

The code allows for the generation of both directed and undirected multi-edge networks with a prescribed strength sequence (optionally including or excluding self-loops) in the various ensembles. In the following section a description of the simulation procedure for each ensemble is given. It must be noted that all node related quantities computed are averaged only over the realizations where a given node exists ($s_i \neq 0$), since the probability for a node to have strength 0 is not zero in the GC and Canonical (C) ensembles (albeit rapidly decreases with \hat{s}).

Micro canonical Ensemble

To generate a multi-edge network in the micro-canonical ensemble, one just needs to apply the well-known configurational model schema [3] (this method was also used in [4]) and *rewire* the link connections permitting multiple connections between nodes. In our case, self-edges have been permitted to allow for an exactly analytical solvable solution of the saddle point equations of the model (see below for a discussion on their effects).

In this case, the expected number of events joining nodes i and j reads,

$$\langle t_{ij} \rangle^{\text{micro}} = \begin{cases} \frac{\hat{s}_i \hat{s}_j}{\hat{T}} & \text{if } i \neq j \\ \frac{\hat{s}_i(\hat{s}_i - 1)}{\hat{T}} & \text{otherwise.} \end{cases} \quad (1)$$

And so for each node, the error committed comparing with the other ensembles is only due to self-loops, $\Delta_{\langle s \rangle} = \frac{\hat{s}_i}{\hat{T}} \leq 1$ which can only be an important quantity for large \hat{s} , but in this case the relative importance of self-loops with respect to the strengths is completely negligible.

It is very important to consider that despite existing many algorithms to *reshuffle* networks preserving the strength sequence, only some of them explore the phase space in an unbiased (or almost unbiased) manner. In our case, the only problem the configurational model has is that it does not allow for self-loops on nodes with strength 1, since they appear always in couples. However, this is not a big issue since the presence of self-loops is only important for high strength nodes while its probability of appearance for nodes with strength 1 is effectively zero.

If one wishes not to accept self-loops, then two alternative approaches are available: one can chose just to not connect the connections corresponding to self-edges (thus not fixing the strength sequence exactly) if rewiring fail repeatedly for sufficiently large number of trials or discard entirely the configurations for which a self-loop event exists.

The first approach induces node correlations in the model and hence the ensemble space is not explored evenly (losing the strict micro canonical nature), it does not fix the strength distribution exactly and furthermore for exponents $\gamma \leq 2$ the number of rejections can grow significantly, making simulations extremely lengthy.

The second approach, on the other hand, is unfeasible in finite time: The probability for a given node of strength \hat{s}_i in the configuration model allowing self-loops while joining stubs to have a self-loop is $p_r^i = \frac{\hat{s}_i - 1 - r}{\hat{T}}$ being r the number of already existing self-loops for the considered node. The number of self-events per node sl_i is thus the result of a random incremental process with step depending probabilities p_r^i . The probability to have no self-loops at the end of the process for each stub is approximately $p_0^i = \left(1 - \frac{\hat{s}_i - 1}{\hat{T}}\right)$ and approximating the states (also self-states) as independent variables, the probability to obtain exactly 0 self loops $SL = \sum_i sl_i$ in the network reads,

$$P(SL = 0) = \prod_i \left(1 - \frac{\hat{s}_i - 1}{\hat{T}}\right)^{\hat{s}_i} \quad (2)$$

(but this is only an upper bound since as stubs get connected, \hat{T} decreases in "time"). In any case, the earlier expression rapidly vanishes due to the large number of terms in the l.h.s. product which are strictly smaller than 1, hence a configuration approach discarding absolutely the configurations with self-edges is not feasible in practice. The complexity of this algorithm is of order $\mathcal{O}(T)$ (the length of the stub sequence). In this scenario, the occupation numbers obtained are, in general, correlated.

Canonical Ensemble

For the canonical ensemble, the statistics of occupation numbers is multinomial with associated probabilities $p_{ij} = \frac{\hat{s}_i \hat{s}_j}{\hat{T}^2}$ and \hat{T} trials. To avoid self-edges one can set $p_{ii} = 0 \forall i$. This method has a limited applicability with system size, since the generation of multinomial distributed variables is not independent and requires a large amount of memory, due to the fact that the occupation numbers generated are correlated,

$$\sigma_{t_{ij}, t_{kl}} = \begin{cases} -Tp_{ij}p_{kl} & ij \neq kl \\ Tp_{ij}(1 - p_{kl}) & ij = kl \end{cases} \quad (3)$$

Grand Canonical Ensemble

The grand canonical ensemble can be implemented in two alternative yet equivalent approaches. One can generate a Poisson distributed number τ with mean T and then generate a collection of occupation numbers $\{t_{ij}\}$ using a multinomial distribution of τ trials and probabilities $p_{ij} = \frac{\hat{s}_i \hat{s}_j}{\hat{T}^2}$ or alternatively one can generate a sequence of L independent Poisson occupation numbers with mean $\langle t_{ij} \rangle = \frac{\hat{s}_i \hat{s}_j}{\hat{T}}$. We have chosen the latter approach to avoid memory overload problems (in this case the occupation numbers are independent and can be generated accordingly). The complexity of this algorithm scales with the number of possible states L . Note that in this case self-edges can be manually avoided by setting $\langle t_{ii} \rangle = 0$ (or equivalently $p_{ii} = 0 \forall i$).

The effect of self-loops

Throughout this letter we have considered ensembles of networks where self-loops are allowed. The reason for this is that in such case the resulting saddle point equations for the hidden variables $\{x_i\}$ can be exactly solved. In the case of not allowing self-loops, then N the saddle point equations (one for each node i) take the form ($X \equiv \sum_j x_j$),

$$\hat{s}_i = \beta x_i \sum_{j \neq i} x_j = \beta x_i (X - x_i), \quad (4)$$

which correspond to a set of non-linear coupled equations and cannot be solved analytically (note that the ensembles are still equivalent if one exactly solves the equations by means of computer methods). If nevertheless we chose to use the solutions for the case of self-loops to this case ($x_i = \hat{s}_i, \beta = \hat{T}^{-1}$), we have that $\langle s \rangle_i = \hat{s}_i \left(1 - \frac{\hat{s}_i}{\hat{T}}\right)$ so the relative error committed is $\epsilon_{\langle s \rangle} = \frac{|\hat{s} - \langle s \rangle|}{\hat{s}} = \frac{\hat{s}}{\hat{T}}$ whose importance depends on the strength of each node but does not vanish in the thermodynamic limit ($T, \hat{s} \rightarrow \infty$).

For non-broad distributed strength sequences with finite mean and standard deviation, this is not a problem even for the worst case scenario, since $\hat{s}_{\max}/\hat{T} \ll 1$ independently of the sampling. Considering the case of skewed distributions (the paradigmatic case power law), we have that the condition $\epsilon_{\max} = \frac{\hat{s}_{\max}}{\hat{T}}$ needs to be analyzed. Re-writing $\hat{T} = N\bar{s}$, and using extreme value theory, we can assess the scaling of this magnitude. For a power law distribution, the maximum value according to the sample can be shown to have a Fréchet distribution,

$$p(\hat{s}_{\max}) = (\gamma - 1) \hat{s}_{\max}^{-\gamma} e^{-\hat{s}_{\max}^{-(\gamma-1)}} \quad \frac{\langle \hat{s}_{\max} \rangle}{N\bar{s}} \sim \frac{N^{\frac{2-\gamma}{\gamma-1}}}{\bar{s}} \quad (5)$$

So we can see that the only problem will come when $\gamma \in (1, 2]$, in which case, we will have to manually set an appropriate value for the ratio $p_{\max} \equiv \hat{s}_{\max}/\hat{T}$ for our calculation to have a good accuracy (bear in mind that since s_{\max} is broadly distributed in such case, the choice will not be random nor general). It is also important to see that in most real cases, very skewed distributions are limited by a natural cut-off that induces an exponential decay on the tail (for physical reasons mainly) so again the condition p_{\max} could be considered mainly negligible.

As a final note, one needs to take into account that the error committed while solving the saddle point equations are of systematic nature. Additionally to this error, one needs to consider the inherent fluctuations of the random variables considered whenever making any sort of calculation of network metrics.

The provided simulation code on the different ensembles allows to accept or not self edges. Please note that all the analytical results provided are approximately valid also for the case of networks without accepting self-loops (given

that the approximation for $p_{max} \ll 1$ early mentioned holds), but one needs to substitute in every expression \hat{s} by $\hat{s} \left(1 - \frac{\hat{s}}{T}\right)$. One also must consider that the sums over indices now do not include the self-connection terms. The main conclusion, however, is that the effect on most network observables will be negligible in general if the value of $p_{max} = \hat{s}_{max}/\hat{T}$ is sufficiently small.

SYNTHETIC POWER LAWS WITH PRESCRIBED AVERAGE STRENGTH

In this paper, to illustrate our findings we have used power law distributed strength sequences. Since in the ensemble approach used here, the thermodynamic limit is defined in terms of $T = \sum s_i$, we use synthetic distributions generated with tunable average strengths. This fact poses problems for the case of power law distributed strengths, since,

$$\bar{s}(\gamma, s_{min}, s_{max}) = \frac{1}{\sum_{s_{min}}^{s_{max}} s^{-\gamma}} \sum_{s_{min}}^{s_{max}} s_i^{1-\gamma}. \quad (6)$$

In such case one has two scenarios: $\gamma > 2$ and $\gamma \leq 2$. For the first scenario, the effect of s_{max} is negligible, and one can achieve an (approximately) desired \bar{s} by setting an appropriate s_{min} . For the second scenario, the opposite happens and one needs to apply a cut-off on s_{max} to limit the average strength of the sequence (which would be unbounded in the case of infinite sampling).

-
- [1] O. Sagarra, C. J. Pérez Vicente, and A. Díaz-Guilera, Phys. Rev. E **88**, 062806 (2013), ISSN 1539-3755, URL <http://link.aps.org/doi/10.1103/PhysRevE.88.062806>.
 - [2] O. Sagarra, F. Font-Clos, C. J. Pérez-Vicente, and a. Díaz-Guilera, EPL (Europhysics Lett. **107**, 38002 (2014), ISSN 0295-5075, URL <http://stacks.iop.org/0295-5075/107/i=3/a=38002?key=crossref.a3c0e51539f4bbc1375fa0f958a695c2>.
 - [3] M. Molloy and B. Reed, Random Struct. Algorithms **6**, 161 (1995), ISSN 10429832, URL <http://doi.wiley.com/10.1002/rsa.3240060204>.
 - [4] M. Á. Serrano, in *AIP Conf. Proc.* (AIP, 2005), vol. 776, pp. 101–107, ISSN 0094243X, URL <http://scitation.aip.org/content/aip/proceeding/aipcp/10.1063/1.1985381>.