

RESEARCH

Open Access



Network reconstruction via density sampling

Tiziano Squartini^{1*} , Giulio Cimini^{1,2}, Andrea Gabrielli^{1,2} and Diego Garlaschelli³

*Correspondence:

tiziano.squartini@imtlucca.it

¹IMT School for Advanced Studies
Lucca, Piazza S.Francesco 19, 55100
Lucca, ItalyFull list of author information is
available at the end of the article

Abstract

Reconstructing weighted networks from partial information is necessary in many important circumstances, e.g. for a correct estimation of systemic risk. It has been shown that, in order to achieve an accurate reconstruction, it is crucial to reliably replicate the empirical degree sequence, which is however unknown in many realistic situations. More recently, it has been found that the knowledge of the degree sequence can be replaced by the knowledge of the strength sequence, which is typically accessible, complemented by that of the total number of links, thus considerably relaxing the observational requirements. Here we further relax these requirements and devise a procedure valid when even the the total number of links is unavailable. We assume that, apart from the heterogeneity induced by the degree sequence itself, the network is homogeneous, so that its (global) link density can be estimated by sampling subsets of nodes with representative density. We show that the best way of sampling nodes is the random selection scheme, any other procedure being biased towards unrealistically large, or small, link densities. We then introduce our core technique for reconstructing both the topology and the link weights of the unknown network in detail. When tested on real economic and financial data sets, our method achieves a remarkable accuracy and is very robust with respect to the sampled subsets, thus representing a reliable practical tool whenever the available topological information is restricted to small portions of nodes.

PACS numbers: 89.75.Hc; 89.65.Gh; 02.50.Tt

Introduction

Reconstructing a weighted, directed network means providing an algorithm to estimate the presence and the weight of all links in the network, making optimal use of the available information (Wells 2004; Upper 2011; Mastromatteo et al. 2012; Baral and Figue 2012; Drehmann and Tarashev 2013; Hałaj and Kok 2013; Anand et al. 2014; Montagna and Lux 2014; Peltonen et al. 2015; Cimini et al. 2015b). Since several networks are in general compatible with the known information, the output of such a procedure cannot identify a unique network but rather an ensemble of possible ones. This leads to a (large) set of candidate networks to be sampled with a certain probability, where the latter has to be specified in such a way that the resulting ensemble average is as close as possible to the empirical, unknown network. Maximum-entropy is a powerful method to construct

probability distributions that realise a certain set of constraints on average. Treating the available pieces of information as empirical constraints in the maximum-entropy procedure ensures that the statistical inference carried out via the resulting distribution is maximally unbiased.

In many situations, e.g. for economic, interbank or other financial networks, the strength sequence (i.e. the list of strengths of all nodes) is known while there is little or no information available about the topology (i.e. the binary structure) of the network. Exploiting the strength sequence as the only constraint of the maximum entropy procedure leads to an unrealistic ensemble where the likely networks are (almost) completely connected (Mastrandrea et al. 2014). This occurs because, when replicating the empirical strengths in absence of topological information, the method tends to distribute non-zero link weights as evenly as possible (i.e. between all pairs of nodes). When such unrealistically dense networks are used as proxies to measure, e.g. the level of systemic risk in a financial network, the resulting estimates are completely unreliable. By contrast, it has been shown that, if the degree sequence is known in addition to the strength sequence, the network reconstruction procedure improves tremendously and achieves a remarkable accuracy, as a result of a much more faithful replication of the topology (Mastrandrea et al. 2014; Cimini et al. 2015a). Notice that, if the link weights are specified by the matrix \mathbf{W} , whose entry $w_{ij} \geq 0$ represents the weight of the directed link from node i to node j , the topology is specified by the binary adjacency matrix \mathbf{A} whose entry $a_{ij} = 1$ if w_{ij} is strictly positive and zero otherwise.

Although complete information on the degree sequence is rarely available, this kind of information can be retrieved from the strength sequence, provided that the latter is complemented with some kind of topological information: in (Musmeci et al. 2013) this information consists of the degree sequence of only a subset I of nodes, $\{k_i\}_{i \in I}$, while in (Cimini et al. 2015b) the information used is the total number of links, L , of the network.

In this paper we face the problem of reconstructing weighted, directed networks, for which the only information available is represented by the set of out-strengths $s_i^{out} = \sum_{j(\neq i)} w_{ij}$ and in-strengths $s_i^{in} = \sum_{j(\neq i)} w_{ji}$ (i.e. the total rows and columns sums of the adjacency matrix) as well as the link density of a subset I of nodes, i.e. $c_I = \frac{L_I}{n_I(n_I-1)}$, with $L_I = \sum_{i \in I} \sum_{j(\neq i) \in I} a_{ij}$ being the observed number of internal links to the subset I . By doing so, we do not require information which is either too detailed (as the degree sequence of even a small subset of nodes) or simply unaccessible (as the total number of links). However, the information encoded into the link density of the chosen subset must be representative of the global one, in order to accurately reconstruct a given network: for this reason, we also propose a recipe about how properly sampling the nodes set of our network. As we will show, the random-nodes sampling scheme provides the best way to draw representative subsets out of the whole nodes set.

Concerning the reconstruction of the weighted structure, we will employ the degree-corrected gravity model (Cimini et al. 2015b) with a correction term ensuring that the strengths are reproduced even in absence of self-loops, i.e. of diagonal terms indicating self-interactions. As we will show, such a correction becomes more and more important as the strength of the considered node is increased, whence the need to properly account for it.

The rest of the paper is organized as follows. In “Methods” section we illustrate the two steps characterizing our reconstruction method and provide measures to test the effectiveness of the algorithm; in “Results” section we apply our method to two real networks, an economic one and a financial one, and in “Conclusions” section we discuss the results.

Methods

Inferring the topological structure

In order to reconstruct the topological structure of a network \mathbf{W} , whenever the nodes strengths $\{s_i^{out}\}_{i=1}^N$ and $\{s_i^{in}\}_{i=1}^N$ and the total number of links L are known, one can follow the algorithm proposed in (Cimini et al. 2015b), which prescribes to solve the equation

$$L = \langle L \rangle \tag{1}$$

with $L = \sum_i \sum_{j(\neq i)} a_{ij}$, $\langle L \rangle = \sum_i \sum_{j(\neq i)} p_{ij}$ and $p_{ij} = (zs_i^{out}s_j^{in})/(1 + zs_i^{out}s_j^{in})$, in order to estimate the unknown parameter z and quantify the probability p_{ij} that a directed link from i to j exists. However, a global (yet very simple) piece of information as L may be not always available. In these cases, an algorithm resorting upon local information has to be employed. In this paper we propose an algorithm to infer the unknown parameter z whenever the information of only a subset I of nodes is accessible. Notice that a possible solution to this problem has already been provided in (Musmeci et al. 2013), where the supposedly known piece of information is represented by the degree sequence of the nodes in I , i.e. $\{k_i\}_{i \in I}$, an hypothesis leading to the equation

$$\sum_{i \in I} (k_i^{out} + k_i^{in}) = \sum_{i \in I} (\langle k_i^{out} \rangle + \langle k_i^{in} \rangle) \tag{2}$$

(with $\langle k_i^{out} \rangle = \sum_{j(\neq i) \in V} p_{ij}$ and $\langle k_i^{in} \rangle = \sum_{j(\neq i) \in V} p_{ji}$ and V indicating the whole nodes set). However, the knowledge of the number of neighbors of even a small subset of nodes may be unavailable as well. For this reason, here we make use of a simpler, more easily accessible, information and suppose to know only the link density within the subset I . Our recipe thus reads

$$c_I = \langle c_I \rangle \tag{3}$$

where $c_I = L_I/[n_I(n_I - 1)]$, $n_I = |I|$ is the number of nodes constituting the subset I , $L_I = \sum_{i \in I} \sum_{j(\neq i) \in I} a_{ij}$ is the observed number of links within it and $\langle L_I \rangle = \sum_{i \in I} \sum_{j(\neq i) \in I} p_{ij}$ is the expected value of L_I .

Remarkably, Eq. (3) can be easily extended to infer the structure of a different subset (say I'), provided that the link density of the latter could be guessed from the known value c_I . As an example, let us assume the existence of a linear proportionality between the two values $c_{I'}$ and c_I : in this case, the equation to be solved would be

$$c_I = f \langle c_{I'} \rangle. \tag{4}$$

More explicitly, such a condition translates into the equation

$$c_I = \frac{f}{n_{I'}(n_{I'} - 1)} \sum_{i \in I'} \sum_{j(\neq i) \in I'} \frac{z_{I'} s_i^{out} s_j^{in}}{1 + z_{I'} s_i^{out} s_j^{in}} \tag{5}$$

which shows that the observed quantity tuning the parameter $z_{I'}$ is $c_I \cdot n_{I'}(n_{I'} - 1)$, i.e. the link density of the known subset, corrected by a volume term.

The value $f = 1$ corresponds to the assumption that the network is homogeneous. This is equivalent to requiring that any two different subsets have exactly the same link density and that, in turn, any subset provides a representative value of the global network density. As we will show in what follows, a random sampling of the set of nodes indeed ensures that this assumption is verified with high accuracy, for the networks considered here.

Inferring the weighted structure

Beside reconstructing a network topological features, the approach proposed in (Cimini et al. 2015b) satisfactorily reproduces also its weighted structure. This approach is based on the degree-corrected gravity model prescription, which reads

$$w_{ij} = \begin{cases} 0 & \text{with probability } 1 - p_{ij}, \\ \frac{s_i^{out} s_j^{in}}{W p_{ij}} & \text{with probability } p_{ij} \end{cases} \tag{6}$$

leading to the expectations $\langle w_{ij} \rangle = s_i^{out} s_j^{in} / W$ and ensuring that $s_i^{out} = \langle s_i^{out} \rangle = \sum_j w_{ij}, \forall i$ and $s_i^{in} = \langle s_i^{in} \rangle = \sum_j w_{ji}, \forall i$ (i.e. that the in-strength and out-strength sequences are, on average, reproduced) as long as *all* entries are summed over.

However, in many real-world networks self-loops are either absent or explicitly excluded: this implies that either the diagonal terms of the adjacency matrix are equal to zero or that our sums should run over $j \neq i$. This causes the expectations coming from the degree-corrected gravity model to need an extra-term to restore the correct value. More explicitly,

$$\langle s_i^{out} \rangle = \sum_{j(\neq i)} \langle w_{ij} \rangle = \frac{s_i^{out} (W - s_i^{in})}{W} = s_i^{out} - \frac{s_i^{out} s_i^{in}}{W}, \tag{7}$$

$$\langle s_i^{in} \rangle = \sum_{j(\neq i)} \langle w_{ji} \rangle = \frac{s_i^{in} (W - s_i^{out})}{W} = s_i^{in} - \frac{s_i^{out} s_i^{in}}{W} \tag{8}$$

and the missing term to be added up to our expectations is precisely the diagonal term, i.e. $\langle w_{ii} \rangle$.

Here we provide a solution to the problem above, by redistributing the diagonal term $\langle w_{ii} \rangle$ across the $N - 1$ entries of the i th row and the $N - 1$ entries of the i th column. In order to implement it, a procedure inspired to the iterative proportional fitting (IPF) algorithm (Bishop et al. 2007) can be devised. More specifically, redistributing the diagonal terms across the corresponding rows and columns amounts to redistribute the strengths of the following matrix on the entries equal to 1. Notice that we need to explicitly distinguish the strengths along rows and columns, since the generic weight w_{ij} needs a correction affecting both i and j .

$$\begin{array}{cccc|c}
 0 & 1 & 1 & 1 & \dots & \frac{s_1^{out} s_1^{in}}{W} \\
 1 & 0 & 1 & 1 & \dots & \frac{s_2^{out} s_2^{in}}{W} \\
 1 & 1 & 0 & 1 & \dots & \frac{s_3^{out} s_3^{in}}{W} \\
 1 & 1 & 1 & 0 & \dots & \frac{s_4^{out} s_4^{in}}{W} \\
 \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\
 \hline
 \frac{s_1^{out} s_1^{in}}{W} & \frac{s_2^{out} s_2^{in}}{W} & \frac{s_3^{out} s_3^{in}}{W} & \frac{s_4^{out} s_4^{in}}{W} & \dots &
 \end{array} \tag{9}$$

In order to achieve the aforementioned redistribution, one can compute the iterations of the IPF algorithm

$$\begin{cases} w_{ij}^{(n)} = \frac{s_i^{out} s_i^{in}}{W} \left(\frac{w_{ij}^{(n-1)}}{\sum_{k(\neq i)} w_{ik}^{(n-1)}} \right) \\ w_{ij}^{(n+1)} = \frac{s_j^{out} s_j^{in}}{W} \left(\frac{w_{ij}^{(n)}}{\sum_{k(\neq j)} w_{kj}^{(n)}} \right) \end{cases} \quad (10)$$

upon setting the matrix defined by $w_{ij}^{(0)} = 1, \forall i \neq j$ as the initial configuration. As a consequence, we need to correct our probabilistic recipe as

$$w_{ij} = \begin{cases} 0 & \text{with probability } 1 - p_{ij}, \\ \left(\frac{s_i^{out} s_j^{in}}{W} + w_{ij}^{(\infty)} \right) \frac{1}{p_{ij}} & \text{with probability } p_{ij}. \end{cases} \quad (11)$$

For all practical purposes, a small number of iterations is often enough to achieve a satisfactory degree of accuracy. Here we explicitly report the analytical functional form of the first three IPF algorithm iterations only:

$$\begin{aligned} w_{ij}^{(1)} &= \frac{s_i^{out} s_i^{in}}{W} \left[\frac{1}{N-1} \right]; \\ w_{ij}^{(2)} &= \frac{s_i^{out} s_i^{in}}{W} \left[\frac{s_j^{out} s_j^{in}}{\sum_{l(\neq j)} s_l^{out} s_l^{in}} \right]; \\ w_{ij}^{(3)} &= \frac{s_i^{out} s_i^{in}}{W} \left[\frac{s_j^{out} s_j^{in}}{\sum_{l(\neq j)} s_l^{out} s_l^{in}} \right] \left[\frac{1}{\sum_{k(\neq i)} \frac{s_k^{out} s_k^{in}}{\sum_{m(\neq k)} s_m^{out} s_m^{in}}} \right]. \end{aligned} \quad (12)$$

A pseudo-code summarizing the two main steps of our algorithm (i.e. Eqs. (3) and (11)) is provided in Appendix 1.

Testing our reconstruction algorithm

An algorithm aiming at reconstructing the topological structure of a network is an example of a binary classifier which tries to infer whether each link is present or not. In order to test the performance of our reconstruction method we, thus, consider four indicators: the number of *true positives*, *true negatives*, *false positives* and *false negatives*. In network terms, the expectation value of such indices reads $\langle TP \rangle = \sum_i \sum_{j(\neq i)} a_{ij} p_{ij}$, $\langle TN \rangle = \sum_i \sum_{j(\neq i)} (1 - a_{ij})(1 - p_{ij})$, $\langle FP \rangle = \sum_i \sum_{j(\neq i)} (1 - a_{ij}) p_{ij}$ and $\langle FN \rangle = \sum_i \sum_{j(\neq i)} a_{ij} (1 - p_{ij})$. However, the information provided by these indicators is often condensed into four alternative indices. The first one is called *sensitivity* (or *true positive rate*), $\langle TPR \rangle = \frac{\langle TP \rangle}{L}$, and quantifies the percentage of 1s that are correctly recovered by our method. The second index is the *specificity* (or *true negative rate*), $\langle SPC \rangle = \frac{\langle TN \rangle}{N(N-1) - L}$, and quantifies the percentage of 0s that are correctly recovered by our method. The third index is the *precision* (or *positive predicted value*), $\langle PPV \rangle = \frac{\langle TP \rangle}{L}$, and measures the performance of our method in correctly placing the 1s with respect to the total number of predicted 1s. The fourth index is the *accuracy*, $\langle ACC \rangle = \frac{\langle TP \rangle + \langle TN \rangle}{N(N-1)}$, and quantifies the overall performance of our method in correctly placing both the 1s and the 0s.

To test the effectiveness of the weighted reconstruction, instead, we use the cosine similarity measure which estimates the distance between the observed weights $\{w_{ij}\}_{i,j=1}^N$ and the conditional expected weights under our model $\{(w_{ij}|a_{ij} = 1)\}_{i,j=1}^N$ by treating

the corresponding matrices as vectors of real numbers and measuring their overlap. In formulas,

$$\theta = \frac{\mathbf{W} \cdot \langle \mathbf{W} \rangle}{\|\mathbf{W}\| \|\langle \mathbf{W} \rangle\|} \quad (13)$$

with $\theta = -1$ indicating maximum dissimilarity, $\theta = 0$ indicating absence of correlations and $\theta = 1$ indicating perfect overlap.

Results

World Trade Web

The first network we have analyzed is the World Trade Web (WTW), i.e. the network whose nodes are the world countries and whose links represent the trade volumes between them: in other words, w_{ij} quantifies the volume of export from i to j . We remand the reader to (Gleditsch 2002) for more details on the dataset. For the sake of illustration, we show detailed results for the snapshot of the WTW in year 2000. We have however analyzed other temporal snapshots as well and found comparable results (see Appendix 2).

Table 1 sums up the results of our analysis when the nodes subset I is chosen at random. We see that the performance of our algorithm is not affected by the cardinality of I upon which the estimation of z is carried on, providing remarkably good results for all the chosen values. In particular, our method is overall very accurate, being able to correctly recover the 80% of 1s and the 73% of 0s, a result to be compared with the performance of a perfect classifier, for which $\langle TPR \rangle = \langle SPC \rangle = 1$, and with that of a random classifier, for which $\langle TPR \rangle = 1 - \langle SPC \rangle = c$ (c being the link density of the whole network). The high accuracy of our reconstruction method is also witnessed by the low rate of false positives of our algorithm, due to the accurate estimation of the actual link density. As discussed in (Squartini et al. 2016), overestimating the link density would have increased the expected TPR (a method predicting a complete network is characterized by $\langle TPR \rangle = 1$), at the price of increasing the rate of false positives as well, thus decreasing the predictive power of the method itself.

Our method performs well also in reproducing the weighted structure of the WTW: upon adding the correction term up to the third iteration of the IPF algorithm, the largest expected in-strength (reading $\langle s_{i,corr}^{in} \rangle = \sum_{j(\neq i)} \left(\frac{s_j^{out} s_i^{in}}{W} + w_{ji}^{(3)} \right)$, $\forall i$) accounts for the 95% of the observed value. On the other hand, the non-corrected value $\langle s_i^{in} \rangle = \sum_{j(\neq i)} \left(\frac{s_j^{out} s_i^{in}}{W} \right)$ accounts for the 82% only. Better results are obtained for the out-strength sequence: the corrected value for the node characterized by the maximum out-strength amounts at the 99% of the corresponding observed value (the non-corrected value accounts for the 88%).

Overall, we obtain a value $\theta_{WTW} \simeq 0.712$ for all the considered cardinalities n_I , indicating a satisfactorily high level of similarity between our weights prediction and their observed values.

e-MID interbank network

The second network we have tested our method upon is the electronic Market for Interbank Deposits (e-MID), i.e. the network whose nodes are banks and whose generic link $i \rightarrow j$ represents the loan granted from i to j . We remand the reader to (Iori et al. 2006) for more details on the dataset.

Table 1 Statistical indicators used to evaluate the performance of our sampled-based reconstruction method, for different cardinalities n of the known subset I . Results are shown together with the 95% confidence intervals (not shown whenever their difference affects the significant digits beyond the third one)

WTW	$n = 5$ (CI 95%)	$n = 10$ (CI 95%)	$n = 20$ (CI 95%)	$n = 50$ (CI 95%)	$n = 100$ (CI 95%)
2000 - True positive rate	0.794 [0.772,0.816]	0.779 [0.765,0.793]	0.804 [0.796,0.812]	0.801 [0.797,0.806]	0.801 [0.799,0.804]
2000 - Specificity	0.700 [0.669,0.731]	0.742 [0.726,0.758]	0.721 [0.710,0.731]	0.728 [0.723,0.734]	0.729 [0.726,0.733]
2000 - Positive predicted value	0.796 [0.784,0.808]	0.810 [0.803,0.817]	0.799 [0.795,0.804]	0.802 [0.800,0.805]	0.802 [0.801,0.803]
2000 - Accuracy	0.755 [0.750,0.760]	0.763 [0.762,0.766]	0.769 [0.768,0.770]	0.771	0.771
2000 - Cosine similarity	0.712	0.712	0.712	0.712	0.712
e-MID					
1999 - True positive rate	0.641 [0.601,0.673]	0.633 [0.614,0.653]	0.633 [0.620,0.646]	0.637 [0.623,0.643]	0.636 [0.632,0.640]
1999 - Specificity	0.839 [0.823,0.856]	0.856 [0.848,0.864]	0.860 [0.854,0.865]	0.860 [0.857,0.863]	0.861 [0.859,0.862]
1999 - Positive predicted value	0.623 [0.611,0.637]	0.632 [0.625,0.639]	0.633 [0.628,0.638]	0.632 [0.623,0.635]	0.633 [0.631,0.634]
1999 - Accuracy	0.785 [0.780,0.790]	0.795 [0.794,0.796]	0.798 [0.797,0.799]	0.799 [0.798,0.800]	0.799
1999 - Cosine similarity	0.810 [0.805,0.815]	0.814 [0.811,0.816]	0.816 [0.815,0.817]	0.817	0.817

The considered cardinalities $n = 5, 10, 20, 50, 100$ correspond to percentages ranging from $\approx 2\%$ to $\approx 50\%$ of the total number of nodes. As reference values, the link density is $c = 0.578$ for the WTW (in the year 2000) and $c = 0.274$ for e-MID (in the year 1999)

Table 1 summarizes the results of our analysis on e-MID in the year 1999 only (again, similar results hold for the other years in our data set - see Appendix). As for the WTW, the performance of our algorithm is not affected by n_I providing again very good results for the whole range of values of the subsets cardinality. In particular, our method is again very accurate, being able to correctly recover the $\simeq 64\%$ of 1s and the $\simeq 86\%$ of 0s. Even if the predictive power of our method is lower than for the WTW case, the accuracy values are comparable, amounting at $\simeq 80\%$.

Our method performs also very well in reproducing the e-MID weighted structure: the correction term coming from the IPF algorithm and calculated for the maximum $\langle s_{i,corr}^{out} \rangle = \sum_{j(\neq i)} \left(\frac{s_i^{out} s_j^{in}}{W} + w_{ij}^{(3)} \right)$, $\forall i$ accounts for the 99% of the observed value. On the other hand, the usual value $\langle s_i^{out} \rangle = \sum_{j(\neq i)} \left(\frac{s_i^{out} s_j^{in}}{W} \right)$ accounts for the 88% only. A comparable result is obtained for the in-strength sequence: the corrected value for the node characterized by the maximum in-strength still amounts at the 99% of the corresponding observed value (the non-corrected value accounts for the 96%).

The value $\theta_{e-MID} \simeq 0.82$ indicates that, on average, a very high level of similarity between observed and predicted weights is again obtained, confirming the degree-corrected gravity model as a good predictor of the links weights.

Random-nodes sampling scheme

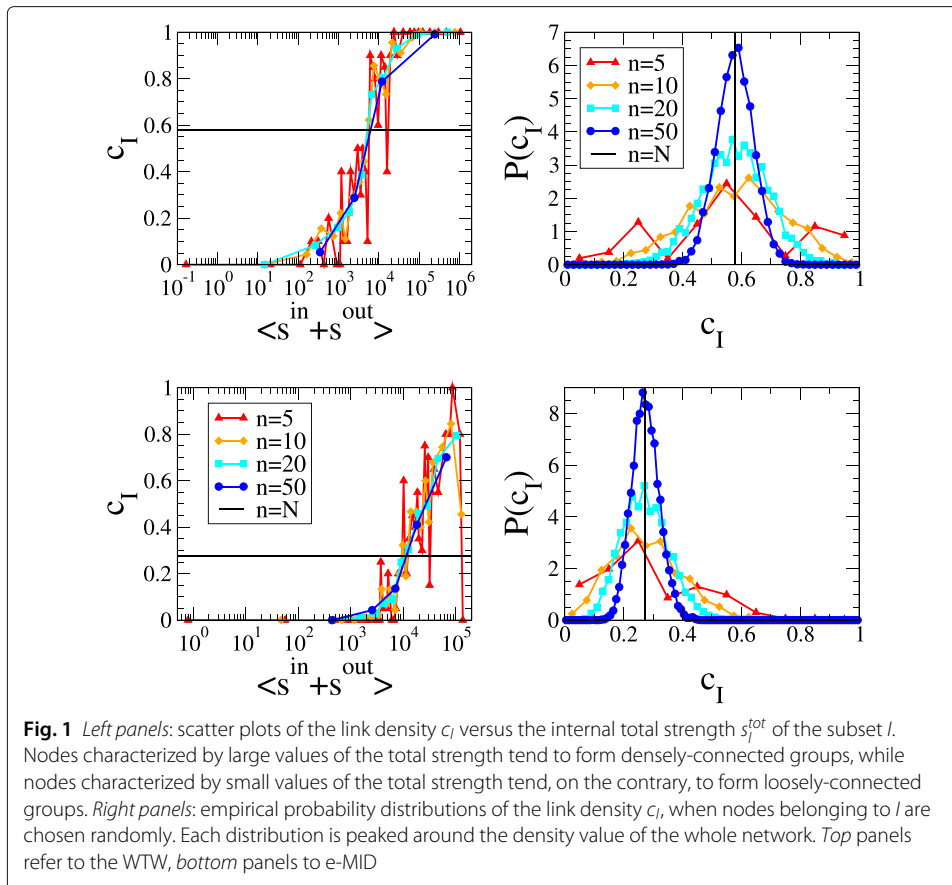
The sampling-based reconstruction algorithm we have proposed in the present paper rests upon the homogeneity assumption, according to which any subset of nodes picked at random provides a representative value of the density of the whole network. Table 2 collects the estimations of the link density, averaged over all sampled subsets of a given cardinality: remarkably, the obtained values are accurate even for low cardinalities. In order to assess the magnitude of fluctuations, we have also explicitly computed the empirical probability distributions of the link density estimates, obtained by random sampling our nodes subsets. These distributions are shown in Fig. 1 (right panels). Naturally, the smaller the cardinality of the considered nodes subsets, the more spaced the values of the observable link density and the less smooth the corresponding probability distribution. These findings suggests that our homogeneity assumption is indeed verified, provided that nodes are sampled according to the random selection scheme (Genois et al. 2015).

As a comparison, we have also sampled nodes sequentially, i.e. by, first, ordering nodes according to their total strength $s_i^{tot} = s_i^{out} + s_i^{in}$ and, then, considering bunches of n subsequent nodes (again, for each value of n). For each subset of nodes we have calculated the corresponding internal link density and plotted it versus the total internal strength of nodes, i.e. $s_i^{tot} = \sum_{i \in I} (s_i^{out} + s_i^{in})$. As shown in Fig. 1 (left panels), such a procedure provides insights on the structural organization of both WTW and e-MID: nodes characterized by large values of the total strength tend to form densely-connected groups whereas nodes characterized by small values of the total strength tend to form loosely-connected groups. Such an evidence confirms the presence of a core-periphery structure, with nodes having a smaller total strength establishing connections with nodes having a large total strength which, in turn, tend to connect preferentially with each other (as a sort of “rich-club”) (Fagiolo et al. 2010; De Masi et al. 2006). Our analysis suggests that a sampling-based reconstruction procedure must rest upon a “balanced” sampling of the nodes, biased neither towards the “core” portion of nodes (which would lead to

Table 2 Link density estimation for different cardinalities n of the random sampled subset A . Results are based on 1000 samples and are shown together with the 95% confidence intervals

Link density	$n = 5$ [CI 95%]	$n = 10$ [CI 95%]	$n = 20$ [CI 95%]	$n = 50$ [CI 95%]	$n = 100$ [CI 95%]
WTW 2000 (true: 0.578)	0.586 [0.560;0.611]	0.559 [0.544;0.574]	0.583 [0.574;0.592]	0.578 [0.573;0.583]	0.577 [0.574;0.580]
e-MID 1999 (true: 0.274)	0.292 [0.271;0.313]	0.278 [0.267;0.289]	0.276 [0.268;0.283]	0.276 [0.272;0.280]	0.275 [0.273;0.278]

The considered cardinalities $n = 5, 10, 20, 50, 100$ correspond to percentages ranging from $\approx 2\%$ to $\approx 50\%$ of the total number of nodes. The true link densities calculated on the entire networks are shown in brackets for reference



severely overestimate the overall network density), nor towards the “periphery” portion of nodes (which would lead to severely underestimate the overall network density). Interestingly, in a recent paper comparing several network sampling techniques was found that the least biased sampling scheme for estimating a given network density is precisely the random-nodes one (Blagus et al. 2016).

Conclusions

The present contribution proposes a recipe to reconstruct a network from a very limited amount of information. In particular, we address the problem of inferring the binary and the weighted structure of a given network from the knowledge of the nodes strengths and the link density of only a subset of nodes. As we have shown in the paper, the best sampling scheme is the random-nodes selection scheme which ensures that an accurate estimation of the whole network density can indeed be achieved. On the contrary, selecting nodes on the basis of more informative structural properties (as the degree, or the strength) could bias the estimation of the connectance towards unrealistically too large, or too small, values. The role played by the available piece of topological information is fundamental not only to achieve an accurate reconstruction of the purely binary structure but also of the weighted structure, as evident upon inspecting Table 1.

The aforementioned results have been obtained by estimating the link density of the whole network upon considering only nodes subsets: in other words, we have verified that

different random subsets (even with different cardinality) are characterized by very similar densities, in turn implying that the whole network density can be estimated (with a high degree of accuracy) by considering a subset randomly drawn from the whole set of nodes. However, the proposed algorithm can be also used to reconstruct networks with a modular structure, upon tuning the link densities of the different modules via Eq. (4): examples are provided by interbank networks structured into jurisdictions, the latter playing the role of the subsets to be reconstructed.

Appendix 1

A pseudo-code summarizing the main steps of the reconstruction algorithm presented in the paper follows.

Algorithm 1: Network reconstruction via density sampling

Input: in- and out-strengths $\{s_i^{in}\}_{i=1}^N$, $\{s_i^{out}\}_{i=1}^N$ and link density of a subset I ,

$$c_I = \frac{L_I}{n_I(n_I-1)}.$$

begin

 solve the equation $c_I = \langle c_I \rangle$ in order to determine z :

$$c_I = \frac{1}{n_I(n_I-1)} \sum_{i \in I} \sum_{j(\neq i) \in I} \frac{z s_i^{out} s_j^{in}}{1 + z s_i^{out} s_j^{in}};$$

 define $p_{ij} = \frac{z s_i^{out} s_j^{in}}{1 + z s_i^{out} s_j^{in}};$

for $m = 1 \dots M$ **do**

for $i < j$ **do**

 calculate the correction to the gravity-like estimation

$$w_{ij}^{(3)} = \frac{s_i^{out} s_j^{in}}{W} \left[\frac{s_j^{out} s_j^{in}}{\sum_{l(\neq j)} s_l^{out} s_l^{in}} \right] \left[\frac{1}{\sum_{k(\neq i)} \frac{s_k^{out} s_k^{in}}{\sum_{m(\neq k)} s_m^{out} s_m^{in}}} \right];$$

 connect i and j with a weight drawn from the following Bernoulli distribution

$$w_{ij} = \begin{cases} 0, & 1 - p_{ij}, \\ \left(\frac{s_i^{out} s_j^{in}}{W} + w_{ij}^{(3)} \right) \frac{1}{p_{ij}}, & p_{ij}. \end{cases}$$

end

end

 verify the goodness of the achieved reconstruction by calculating the ensemble average of indicators like TPR , SPC , PPV , ACC and θ ;

end

Output: ensemble of M reconstructed directed, weighted networks.

Appendix 2

Additional years have been analysed for both the WTW and e-MID (see Tables 3 and 4).

Table 3 Statistical indicators used to evaluate the performance of our sampled-based reconstruction method, for different cardinalities n of the known subset I . Results are shown together with the 95% confidence intervals (not shown whenever their difference affects the significant digits beyond the third one)

WTW	$n = 5$ (CI 95%)	$n = 10$ (CI 95%)	$n = 20$ (CI 95%)	$n = 50$ (CI 95%)
1950 - Link density (true: 0.402)	0.401 [0.375;0.426]	0.402 [0.387;0.416]	0.401 [0.393;0.409]	0.400 [0.396;0.403]
1950 - Accuracy	0.736 [0.731;0.741]	0.747 [0.746;0.749]	0.751	0.752
1950 - Cosine similarity	0.460 [0.458;0.462]	0.463	0.463	0.463
1960 - Link density (true: 0.383)	0.329 [0.305;0.352]	0.343 [0.330;0.357]	0.346 [0.338;0.355]	0.348 [0.344;0.353]
1960 - Accuracy	0.737 [0.734;0.741]	0.746	0.749 [0.748;0.750]	0.751
1960 - Cosine similarity	0.586	0.591	0.591	0.591
1970 - Link density (true: 0.460)	0.464 [0.436;0.492]	0.478 [0.462;0.496]	0.461 [0.451;0.471]	0.464 [0.458;0.469]
1970 - Accuracy	0.695 [0.691;0.699]	0.704 [0.702;0.706]	0.709	0.709
1970 - Cosine similarity	0.669	0.669	0.669	0.669
1980 - Link density (true: 0.468)	0.484 [0.458;0.510]	0.470 [0.455;0.485]	0.471 [0.461;0.481]	0.463 [0.458;0.469]
1980 - Accuracy	0.719 [0.715;0.723]	0.731 [0.730;0.733]	0.734 [0.733;0.735]	0.736
1980 - Cosine similarity	0.732	0.732	0.732	0.732
1990 - Link density (true: 0.505)	0.495 [0.467;0.522]	0.516 [0.500;0.532]	0.506 [0.497;0.515]	0.507 [0.503;0.512]
1990 - Accuracy	0.731 [0.726;0.736]	0.743 [0.741;0.745]	0.748	0.749
1990 - Cosine similarity	0.751	0.751	0.751	0.751

The considered cardinalities $n = 5, 10, 20, 50$ correspond to percentages ranging from $\simeq 2\%$ to $\simeq 25\%$ of the total number of nodes. The true link densities calculated on the entire networks for the various periods are shown in brackets for reference

Table 4 Statistical indicators used to evaluate the performance of our sampled-based reconstruction method, for different cardinalities n of the known subset I . Results are shown together with the 95% confidence intervals (not shown whenever their difference affects the significant digits beyond the third one)

e-MID	$n = 5$ (CI 95%)	$n = 10$ (CI 95%)	$n = 20$ (CI 95%)	$n = 50$ (CI 95%)
2000 - Link density (true: 0.278)	0.293 [0.269;0.317]	0.279 [0.263;0.295]	0.273 [0.264;0.281]	0.280 [0.273;0.283]
2000 - Accuracy	0.763 [0.759;0.768]	0.772 [0.769;0.775]	0.778 [0.777;0.779]	0.778 [0.777;0.779]
2000 - Cosine similarity	0.573 [0.566;0.580]	0.578 [0.576;0.582]	0.582	0.582
2001 - Link density (true: 0.263)	0.279 [0.256;0.303]	0.264 [0.249;0.278]	0.257 [0.246;0.267]	0.266 [0.261;0.272]
2001 - Accuracy	0.763 [0.757;0.770]	0.774 [0.772;0.776]	0.777 [0.775;0.779]	0.778 [0.777;0.779]
2001 - Cosine similarity	0.560 [0.554;0.566]	0.566 [0.563;0.569]	0.569	0.570
2002 - Link density (true: 0.233)	0.253 [0.230;0.276]	0.237 [0.221;0.252]	0.235 [0.225;0.246]	0.233 [0.228;0.239]
2002 - Accuracy	0.759 [0.752;0.766]	0.767 [0.763;0.771]	0.770 [0.767;0.772]	0.772 [0.770;0.773]
2002 - Cosine similarity	0.684 [0.675;0.694]	0.670 [0.682;0.697]	0.699 [0.697;0.701]	0.701 [0.700;0.702]
2003 - Link density (true: 0.214)	0.248 [0.223;0.273]	0.225 [0.208;0.243]	0.217 [0.205;0.228]	0.213 [0.208;0.219]
2003 - Accuracy	0.746 [0.737;0.756]	0.758 [0.752;0.763]	0.763 [0.759;0.766]	0.766 [0.764;0.767]
2003 - Cosine similarity	0.461 [0.453;0.470]	0.462 [0.454;0.470]	0.472 [0.469;0.475]	0.476 [0.475;0.477]
2004 - Link density (true: 0.190)	0.210 [0.185;0.235]	0.183 [0.168;0.199]	0.194 [0.182;0.205]	0.187 [0.181;0.192]
2004 - Accuracy	0.772 [0.762;0.783]	0.785 [0.780;0.790]	0.784 [0.780;0.788]	0.788 [0.786;0.790]
2004 - Cosine similarity	0.481 [0.470;0.492]	0.482 [0.474;0.491]	0.497 [0.493;0.502]	0.503 [0.501;0.504]
2005 - Link density (true: 0.201)	0.232 [0.205;0.258]	0.210 [0.190;0.222]	0.210 [0.200;0.221]	0.208 [0.203;0.214]
2005 - Accuracy	0.751 [0.740;0.762]	0.767 [0.760;0.773]	0.767 [0.763;0.771]	0.769 [0.767;0.771]
2005 - Cosine similarity	0.461 [0.448;0.474]	0.476 [0.470;0.483]	0.486 [0.483;0.490]	0.491 [0.490;0.492]

The considered cardinalities $n = 5, 10, 20, 50, 100$ correspond to percentages ranging from $\simeq 2\%$ to $\simeq 25\%$ of the total number of nodes. The true link densities calculated on the entire networks for the various periods are shown in brackets for reference

Acknowledgments

This work was supported by the EU projects CoeGSS (grant num. 676547), Multiplex (grant num. 317532), Shakermaker (grant num. 687941), SoBigData (grant num. 654024), GrowthCom (grant num. 611272) and the FET projects SIMPOL (grant num. 610704) and DOLFINS (grant num. 640772). AG acknowledges the CNR Strategic Project CRISISLAB funded by Italian Government. DG acknowledges support from the Econophysics foundation (Stichting Econophysics, Leiden, the Netherlands).

Authors' contributions

TS, GC, AG and DG participated in the design of the analysis. TS and GC performed the statistical analysis. All authors wrote, read and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Author details

¹IMT School for Advanced Studies Lucca, Piazza S.Francesco 19, 55100 Lucca, Italy. ²Istituto dei Sistemi Complessi (ISC) - CNR, UoS Sapienza, Dipartimento di Fisica, Università "Sapienza", Piazzale Aldo Moro 5, 00185 Rome, Italy.

³Instituut-Lorentz for Theoretical Physics, Leiden Institute of Physics, University of Leiden, Niels Bohrweg 2, 2333 CA Leiden, The Netherlands.

Received: 20 October 2016 Accepted: 11 January 2017

Published online: 28 January 2017

References

- Anand K, Craig B, von Peter G (2014) Filling in the blanks: Network structure and interbank contagion. *Quant Finance* 15:625–636
- Baral P, Figue JP (2012) Estimation of bilateral exposures - a copula approach. Mimeo. http://www.cirano.qc.ca/conferences/public/pdf/networks2012/02-BARAL-FIQUE-Estimation_of_Bilateral_Exposures-A_Copula_Approach.pdf
- Bishop YM, Fienberg SE, Holland PW (2007) *Discrete Multivariate Analysis: Theory and Practice*. Springer-Verlag, New York
- Blagus N, Subelj L, Bajec M (2016) Empirical comparison of network sampling techniques. arxiv:1506.02449. <https://arxiv.org/pdf/1506.02449v2.pdf>
- Cimini G, Squartini T, Gabrielli A, Garlaschelli D (2015a) Estimating topological properties of weighted networks from limited information. *Phys Rev E* 92:040802. doi:10.1103/PhysRevE.92.040802
- Cimini G, Squartini T, Garlaschelli D, Gabrielli A (2015b) Systemic risk analysis on reconstructed economic and financial networks. *Sci Rep* 5:15758. doi:10.1038/srep15758
- De Masi G, Iori G, Caldarelli G (2006) Fitness model for the Italian interbank money market. *Phys Rev E* 74:066112. doi:10.1103/PhysRevE.74.066112
- Drehmann M, Tarashev N (2013) Measuring the systemic importance of interconnected banks. *J Financ Intermediation* 22(4):586–607. doi:10.1016/j.jfi.2013.08.001
- Fagiolo G, Reyes J, Schiavo S (2010) The evolution of the world trade web: a weighted-network analysis. *J Evol Econ* 20(4):479–514. doi:10.1007/s00191-009-0160-x
- Genois M, Vestergaard C, Cattuto C, Barrat A (2015) Compensating for population sampling in simulations of epidemic spread on temporal contact networks. *Nat Commun* 6(8860)
- Gleditsch KS (2002) Expanded trade and GDP data. *J Confl Resol* 46(5):712–724
- Hałaj G, Kok C (2013) Assessing interbank contagion using simulated networks. *Comput Manag Sci* 10(2):157–186. doi:10.1007/s10287-013-0168-4
- Iori G, Jafarey S, Padilla FG (2006) Systemic risk on the interbank market. *J Econ Behav Organ* 61(4):525–542. doi:10.1016/j.jebo.2004.07.018
- Mastrandrea R, Squartini T, Fagiolo G, Garlaschelli D (2014) Enhanced reconstruction of weighted networks from strengths and degrees. *New J Phys* 16(4):043022
- Mastromatteo I, Zarinelli E, Marsili M (2012) Reconstruction of financial networks for robust estimation of systemic risk. *J Stat Mech Theory Exp* 03:03011. doi:10.1088/1742-5468/2012/03/P03011
- Montagna M, Lux T (2014) Contagion risk in the interbank market: A probabilistic approach to cope with incomplete structural information. Kiel working papers, Kiel Institute for the World Economy. https://www.ifw-members.ifw-kiel.de/publications/1937-contagion-risk-in-the-interbank-market-a-probabilistic-approach-to-cope-with-incomplete-structural-information/1937_KWP.pdf
- Musmeci N, Battiston S, Caldarelli G, Puliga M, Gabrielli A (2013) Bootstrapping topological properties and systemic risk of complex networks using the fitness model. *J Stat Phys* 151(3):720–734. doi:10.1007/s10955-013-0720-1
- Peltonen TA, Rancan M, Sarlin P (2015) Interconnectedness of the banking sector as a vulnerability to crises. Working Paper 1866, ECB European Central Bank. <http://www.ecb.europa.eu/pub/pdf/scpwps/ecbwp1866.en.pdf>
- Squartini T, Caldarelli G, Cimini G (2016) Stock markets reconstruction via entropy maximization driven by fitness and density. arXiv:1606.07684. <https://arxiv.org/pdf/1606.07684v1.pdf>
- Upper C (2011) Simulation methods to assess the danger of contagion in interbank markets. *J Financ Stability* 7(3):111–125. doi:10.1016/j.jfs.2010.12.001
- Wells SJ (2004) Financial interlinkages in the United Kingdom's interbank market and the risk of contagion. Bank of England Working Paper 230, Bank of England. <http://www.bankofengland.co.uk/archive/Documents/historicpubs/workingpapers/2004/wp230.pdf>