

# Network reconstruction and community detection from dynamics

Tiago P. Peixoto\*

Department of Network and Data Science, Central European University, H-1051 Budapest, Hungary  
ISI Foundation, Via Chisola 5, 10126 Torino, Italy and  
Department of Mathematical Sciences, University of Bath,  
Claverton Down, Bath BA2 7AY, United Kingdom

We present a scalable nonparametric Bayesian method to perform network reconstruction from observed functional behavior that at the same time infers the communities present in the network. We show that the joint reconstruction with community detection has a synergistic effect, where the edge correlations used to inform the existence of communities are also inherently used to improve the accuracy of the reconstruction which, in turn, can better inform the uncovering of communities. We illustrate the use of our method with observations arising from epidemic models and the Ising model, both on synthetic and empirical networks, as well as on data containing only functional information.

The observed functional behavior of a wide variety large-scale system is often the result of a network of pairwise interactions. However, in many cases these interactions are hidden from us, either because they are impossible to measure directly, or because their measurement can be done only at significant experimental cost. Examples include the mechanisms of gene and metabolic regulation [1], brain connectivity [2], the spread of epidemics [3], systemic risk in financial institutions [4], and influence in social media [5]. In such situations, we are required to *infer* the network of interactions from the observed functional behavior. Researchers have approached this reconstruction task from a variety of angles, resulting in many different methods, including thresholding the correlation between time-series [6], inversion of deterministic dynamics [7–9], statistical inference of graphical models [10–14] and of models of epidemic spreading [15–20], as well as approaches that avoid explicit modeling, such as those based on transfer entropy [21], Granger causality [22], compressed sensing [23–25], generalized linearization [26], and matching of pairwise correlations [27, 28].

In this work, we approach the problem of network reconstruction in a manner that is different from the aforementioned methods in two important ways. First, we employ a nonparametric Bayesian formulation of the problem, which yields a full posterior distribution of possible networks that are compatible with the observed dynamical behavior. Second, we perform network reconstruction jointly with *community detection* [29], where at the same time as we infer the edges of the underlying network, we also infer its modular structure [30]. As we will show, while network reconstruction and community detection are desirable goals on their own, joining these two tasks has a synergistic effect, whereby the detection of communities significantly increases the accuracy of the reconstruction, which in turn improves the discovery of the communities, when compared to performing these tasks in isolation.

Some other approaches combine community detection with functional observation. Berthet et al. [31] derived necessary conditions for the exact recovery of group assignments for dense weighted networks generated with community structure given observed microstates of an Ising model. Hoffmann et al. [32] proposed a method to infer community structure from time-series data that bypasses network reconstruction, by employing instead a direct modeling of the dynamics given the group assignments. However, neither of these approaches attempt to perform network reconstruction together with community detection. Furthermore, they are tied down to one particular inverse problem, and as we will show, our general approach can be easily extended to an open-ended variety of functional models.

*Bayesian network reconstruction* — We approach the network reconstruction task similarly to the situation where the network edges are measured directly, but via an uncertain process [33, 34]: If  $\mathcal{D}$  is the measurement of some process that takes place on a network, we can define a posterior distribution for the underlying adjacency matrix  $\mathbf{A}$  via Bayes’ rule,

$$P(\mathbf{A}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{A})P(\mathbf{A})}{P(\mathcal{D})}, \quad (1)$$

where  $P(\mathcal{D}|\mathbf{A})$  is an arbitrary *forward* model for the dynamics given the network,  $P(\mathbf{A})$  is the prior information on the network structure, and  $P(\mathcal{D}) = \sum_{\mathbf{A}} P(\mathcal{D}|\mathbf{A})P(\mathbf{A})$  is a normalization constant comprising the total evidence for the data  $\mathcal{D}$ . We can unite reconstruction with community detection via an, at first, seemingly minor, but ultimately consequential modification of the above equation, where we introduce a structured prior  $P(\mathbf{A}|\mathbf{b})$  where  $\mathbf{b}$  represents the partition of the network in communities, i.e.  $\mathbf{b} = \{b_i\}$ , where  $b_i \in \{1, \dots, B\}$  is group membership of node  $i$ . This partition is unknown, and is inferred together with the network itself, via the joint posterior distribution

$$P(\mathbf{A}, \mathbf{b}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{A})P(\mathbf{A}|\mathbf{b})P(\mathbf{b})}{P(\mathcal{D})}. \quad (2)$$

\* peixoto@ceu.edu

The prior  $P(\mathbf{A}|\mathbf{b})$  is an assumed generative model for the network structure. In our work, we will use the degree-corrected stochastic block model (DC-SBM) [35], which assumes that, besides differences in degree, nodes belonging to the same group have statistically equivalent connection patterns, according to the joint probability

$$P(\mathbf{A}|\boldsymbol{\lambda}, \boldsymbol{\kappa}, \mathbf{b}) = \prod_{i < j} \frac{e^{-\kappa_i \kappa_j \lambda_{b_i, b_j}} (\kappa_i \kappa_j \lambda_{b_i, b_j})^{A_{ij}}}{A_{ij}!}, \quad (3)$$

with  $\lambda_{rs}$  determining the average number of edges between groups  $r$  and  $s$  and  $\kappa_i$  the average degree of node  $i$ . The marginal prior is obtained by integrating over all remaining parameters weighted by their respective prior distributions,

$$P(\mathbf{A}|\mathbf{b}) = \int P(\mathbf{A}|\boldsymbol{\lambda}, \boldsymbol{\kappa}, \mathbf{b}) P(\boldsymbol{\kappa}|\mathbf{b}) P(\boldsymbol{\lambda}|\mathbf{b}) d\boldsymbol{\kappa} d\boldsymbol{\lambda}. \quad (4)$$

which can be computed exactly for standard prior choices, although it can be modified to include hierarchical priors that have an improved explanatory power [36] (see Appendix A for a concise summary).

The use of the DC-SBM as a prior probability in Eq. 2 is motivated by its ability to inform link prediction in networks where some fraction of edges have not been observed or have been observed erroneously [34, 38]. The latent conditional probabilities of edges existing between groups of nodes is learned by the collective observation of many similar edges, and these correlations are leveraged to extrapolate the existence of missing or spurious ones. The same mechanism is expected to aid the reconstruction task, where edges are not observed directly, but the observed functional behavior yields a posterior distribution on them, allowing the same kind of correlations to be used as an additional source of evidence for the reconstruction, going beyond what the dynamics alone says.

Our reconstruction approach is finalized by defining an appropriate model for the functional behavior, determining  $P(\mathcal{D}|\mathbf{A})$ . Here we will consider two kinds of indirect data. The first comes from a SIS epidemic spreading model [39], where  $\sigma_i(t) = 1$  means node  $i$  is infected at time  $t$ , 0 otherwise. The likelihood for this model is

$$P(\boldsymbol{\sigma}|\mathbf{A}, \boldsymbol{\tau}, \gamma) = \prod_t \prod_i P(\sigma_i(t)|\boldsymbol{\sigma}(t-1)), \quad (5)$$

where

$$P(\sigma_i(t)|\boldsymbol{\sigma}(t-1)) = f(e^{m_i(t-1)}, \sigma_i(t))^{1-\sigma_i(t-1)} \times f(\gamma, \sigma_i(t))^{\sigma_i(t-1)} \quad (6)$$

is the transition probability for node  $i$  at time  $t$ , with  $f(p, \sigma) = (1-p)^\sigma p^{1-\sigma}$ , and where  $m_i(t) = \sum_j A_{ij} \ln(1 - \tau_{ij}) \sigma_j(t)$  is the contribution from all neighbors of node  $i$  to its infection probability at time  $t$ . In the equations above the value  $\tau_{ij}$  is the probability of an infection via an existing edge  $(i, j)$ , and  $\gamma$  is the  $1 \rightarrow 0$  recovery probability. With these additional parameters, the full posterior

distribution for the reconstruction becomes

$$P(\mathbf{A}, \mathbf{b}, \boldsymbol{\tau}|\boldsymbol{\sigma}) = \frac{P(\boldsymbol{\sigma}|\mathbf{A}, \boldsymbol{\tau}, \gamma) P(\mathbf{A}|\mathbf{b}) P(\mathbf{b}) P(\boldsymbol{\tau})}{P(\boldsymbol{\sigma}|\gamma)}. \quad (7)$$

Since  $\tau_{ij} \in [0, 1]$  we use the uniform prior  $P(\boldsymbol{\tau}) = 1$ . Note also that the recovery probability  $\gamma$  plays no role on the reconstruction algorithm, since its term in the likelihood does not involve  $\mathbf{A}$  (and hence, gets cancelled out in the denominator  $P(\boldsymbol{\sigma}|\gamma) = P(\gamma|\boldsymbol{\sigma})P(\boldsymbol{\sigma})/P(\gamma)$ ). This means that the above posterior only depends on the infection events  $0 \rightarrow 1$ , and thus is also valid without any modifications to all epidemic variants SI, SIR, SEIR, etc [39], since the infection events occur with the same probability for all these models.

The second functional model we consider is the Ising model, where spin variables on the nodes  $\mathbf{s} \in \{-1, 1\}^N$  are sampled according to the joint distribution

$$P(\mathbf{s}|\mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) = \frac{\exp\left(\beta \sum_{i < j} J_{ij} A_{ij} s_i s_j + \sum_i h_i s_i\right)}{Z(\mathbf{A}, \beta, \mathbf{J}, \mathbf{h})}, \quad (8)$$

where  $\beta$  is the inverse temperature,  $J_{ij}$  is the coupling on edge  $(i, j)$ ,  $h_i$  is a local field on node  $i$ , and  $Z(\mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) = \sum_{\mathbf{s}} \exp(\beta \sum_{i < j} J_{ij} A_{ij} s_i s_j + \sum_i h_i s_i)$  is the partition function. Note that this is not a dynamical model, as each microstate  $\mathbf{s}$  is sampled independently according to the above distribution. Unlike the SIS model considered before, this distribution cannot be written in closed form since  $Z(\mathbf{A}, \beta, \mathbf{J}, \mathbf{h})$  cannot be computed exactly, rendering the reconstruction problem intractable. Therefore, instead, we make use of the pseudolikelihood approximation [40], which is very accurate for the purpose at hand [14], where we approximate Eq. 8 as a product of (properly normalized) conditional probabilities for each spin variable  $s_i$

$$P(\mathbf{s}|\mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) = \prod_i \frac{\exp(\beta s_i \sum_j J_{ij} A_{ij} s_j + h_i s_i)}{2 \cosh(\beta \sum_j J_{ij} A_{ij} s_j + h_i)}. \quad (9)$$

With the above likelihood, reconstruction is performed by observing a set of  $M$  microstates  $\bar{\mathbf{s}} = \{\mathbf{s}_1, \dots, \mathbf{s}_M\}$ , sampled according to  $P(\bar{\mathbf{s}}|\mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) = \prod_t P(\mathbf{s}_t|\mathbf{A}, \beta, \mathbf{J}, \mathbf{h})$ , which yields the posterior distribution

$$P(\mathbf{A}, \mathbf{b}, \beta, \mathbf{J}, \mathbf{h}|\bar{\mathbf{s}}) = \frac{P(\bar{\mathbf{s}}|\mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) P(\beta) P(\mathbf{h}) P(\mathbf{J}|\mathbf{A}) P(\mathbf{A}|\mathbf{b}) P(\mathbf{b})}{P(\bar{\mathbf{s}})}. \quad (10)$$

In the above we use uniform priors  $P(\mathbf{J}|\mathbf{A}) = \prod_{ij} [-1/2 < J_{ij} < 1/2]^{A_{ij}}$ , thus forcing, without loss of generality, the values of  $J_{ij}$  to lie in the shifted unit interval  $[-1/2, 1/2]$ . For the remaining parameters we use uniform priors,  $P(\mathbf{h}) \propto 1$  and  $P(\beta) \propto 1$ , for  $\beta \in [-\infty, \infty]$  and  $\mathbf{h} \in [-\infty, \infty]^N$ .

For any of the above posterior distributions, we perform sampling using Markov chain Monte Carlo

(MCMC): For each proposal  $\mathbf{A} \rightarrow \mathbf{A}'$ , it is accepted with the Metropolis-Hastings probability [41, 42]

$$\min \left( 1, \frac{P(\mathbf{A}', \mathbf{b}, \boldsymbol{\theta} | \mathcal{D}) P(\mathbf{A}' \rightarrow \mathbf{A})}{P(\mathbf{A}, \mathbf{b}, \boldsymbol{\theta} | \mathcal{D}) P(\mathbf{A} \rightarrow \mathbf{A}')} \right)$$

and likewise for the node partition  $\mathbf{b} \rightarrow \mathbf{b}'$ , and any of the remaining parameters  $\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}'$ . Note that the acceptance probability does not require the intractable normalization constant  $P(\mathcal{D})$  to be computed. For both functional models considered, a whole sweep over  $E$  entries of the adjacency matrix and  $N$  nodes is done in time  $O(EM + N\langle k \rangle)$ , where  $M$  is the number of data samples per node, allowing the method to be applied for large systems. We summarize and give more details about the technical aspects of the algorithm in Appendix C.

*Synthetic networks* — We begin by investigating the reconstruction performance of networks sampled from the planted partition model (PP), i.e. a DC-SBM with  $\kappa_i = 1$ ,  $\lambda_{rs} = \lambda_{\text{in}}\delta_{rs} + \lambda_{\text{out}}(1 - \delta_{rs})$ , with  $\lambda_{\text{in}} = \langle k \rangle(1 + \epsilon(B - 1))/N$  and  $\lambda_{\text{out}} = \langle k \rangle(1 - \epsilon)/N$ , where  $\epsilon = N(\lambda_{\text{in}} - \lambda_{\text{out}})/\langle k \rangle B$  controls the strength of the modular structure. The detectability threshold for this model is given by  $\epsilon^* = 1/\sqrt{\langle k \rangle}$ , below which it is impossible to recover the planted community structure [43]. Given a network  $\mathbf{A}^*$  from this model, we sample  $M$  independent Ising microstates  $\mathbf{s}$  according to Eq. 8 using  $J_{ij} = 1$ ,  $h_i = 0$  and  $\beta = \beta^*$  being the critical inverse temperature for the particular network. We compare two inference approaches: In the first we sample both the reconstructed network as well as its community structure form the joint posterior of Eq. 10. In the second approach, we perform reconstruction and community detection separately, by first performing reconstruction in isolation, by replacing the DC-SBM prior  $P(\mathbf{A}|\mathbf{b})$  by the likelihood of an Erdős-Rényi model. We evaluate the quality of the reconstruction via the posterior similarity  $S \in [0, 1]$ , defined as

$$S(\mathbf{A}^*, \boldsymbol{\pi}) = 1 - \frac{\sum_{i < j} |A_{ij}^* - \pi_{ij}|}{\sum_{i < j} |A_{ij}^* + \pi_{ij}|}, \quad (11)$$

where  $\mathbf{A}^*$  is the true network and  $\boldsymbol{\pi}$  is the marginal posterior probability for each edge, i.e.  $\pi_{ij} = \sum_{\mathbf{A}, \mathbf{b}, \boldsymbol{\theta}} A_{ij} P(\mathbf{A}, \mathbf{b}, \boldsymbol{\theta} | \mathcal{D})$ . A value  $S = 1$  means perfect reconstruction. We then perform community detection *a posteriori* by obtaining the maximum marginal point estimate

$$\hat{A}_{ij} = \begin{cases} 1 & \text{if } \pi_{ij} > 1/2, \\ 0 & \text{if } \pi_{ij} < 1/2. \end{cases} \quad (12)$$

and then sampling from the posterior  $P(\mathbf{b} | \hat{\mathbf{A}})$ . Fig. 1 contains the comparison between both approaches for networks sampled from the PP model, which shows how sampling from the joint posterior improves both the reconstruction as well as community detection. For the latter, the joint inference allows the detection all the way down to the detectability threshold, for the examples considered, which, otherwise, is not possible with the separate method.

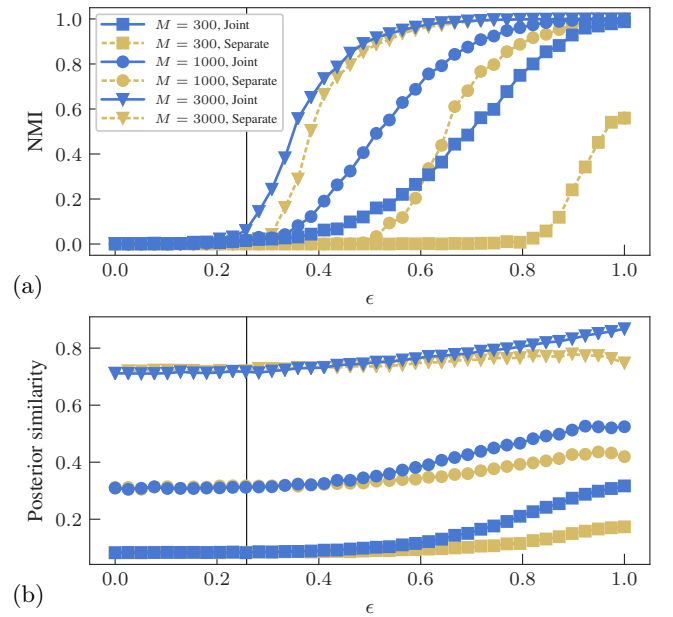


Figure 1. Comparison between joint and separate reconstruction with community detection for a PP model with  $N = 1000$ ,  $\langle k \rangle = 15$  and  $B = 10$ . (a) Normalized mutual information (NMI) between inferred and planted node partitions, as a function of the model parameter  $\epsilon$ , for several values of the number of samples  $M$  from the Ising model described in the text. (b) Posterior similarity between planted and inferred networks, for the same cases as in (a). The vertical line marks the detectability threshold  $\epsilon = 1/\sqrt{\langle k \rangle}$ .

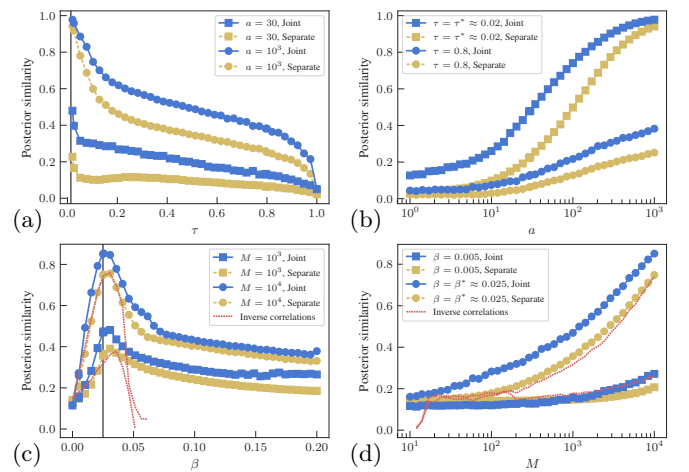


Figure 2. Reconstruction results for simulated dynamics on empirical networks, comparing separate and joint reconstruction with community detection. (a) and (b) correspond to a SIS dynamics on global airport data, using  $\tau_{ij} = \tau$ ,  $\gamma = 1$ , for different values of the infection probability  $\tau$  and node activity  $a$  (defined as the number of infection events per node), and (c) and (d) the Ising model on a food web, using  $J_{ij} = 1$  and  $h_i = 0$ . The dashed red line corresponds to the inverse correlation method for the Ising model. The solid vertical line marks the critical value for each model.

*Real networks with synthetic dynamics* — Now, we investigate the reconstruction of networks not generated by the DC-SBM. We take two empirical networks, the worldwide network of  $N = 3\,286$  airports [44] with  $E = 39\,430$  edges, and a food web from Little Rock Lake [45], containing  $N = 183$  nodes and  $E = 2\,434$  edges, and we sample from the SIS (mimicking the spread of a pandemic) and Ising model (representing simplified inter-species interactions) on them, respectively, and evaluate the reconstruction obtained via the joint and separate inference with community detection, with results shown in Fig. 2. As is also the case for synthetic networks, the reconstruction quality is significantly improved by performing joint community detection [46]. The quality of the reconstruction peaks at the critical threshold for each model, at which the sensitivity to perturbations is the largest. As the number of observed samples increases, so does the quality of the reconstruction, and the relative advantage of the joint reconstruction diminishes, as the data eventually “washes out” the contribution from the prior. For the Ising model, we compare the results of our method with the mean-field inverse correlations method [14], i.e.  $\beta A_{ij} J_{ij} = [C^{-1}]_{ij}$ , where  $C_{ij} = \langle \sigma_i \sigma_j \rangle - \langle \sigma_i \rangle \langle \sigma_j \rangle$  is the connected correlation matrix. As seen in Fig. 2, this simpler reconstruction method can be just as accurate as our separate reconstruction approach, but only close to the critical point. For higher inverse temperatures the reconstruction deteriorates rapidly, and breaks down completely as the system becomes locally magnetized, with whole rows and columns of the matrix  $C$  being equal to zero, causing it to be singular [47]. In such situations this kind of approach requires explicit regularization techniques [48], which become unnecessary with our Bayesian method. The joint inference with community structure improves the reconstruction even further, beyond what is obtainable with typical inverse Ising methods, since it incorporates a different source of evidence.

In Fig. 3 we show a comparison of the reconstruction of the food web network from a simulated Ising model, using different approaches. Optimal thresholding corresponds to the naive approach of imputing the existence of an edge to the connected correlation between two nodes exceeding a threshold  $c^*$ , i.e.  $\pi_{ij} = \{1 \text{ if } C_{ij} > c^*, 0 \text{ otherwise}\}$ . The value of  $c^*$  was chosen to maximize the posterior similarity, which represents the best possible reconstruction achievable with this method. Nevertheless, the network thus obtained is severely distorted. The inverse correlation method comes much closer to the true network, but is superseded by the joint inference with community detection.

*Empirical dynamics* — We turn to the reconstruction from observed empirical dynamics with unknown underlying interactions. The first example is the sequence of  $M = 619$  votes of  $N = 575$  deputies in the 2007 to 2011 session of the lower chamber of the Brazilian congress. Each deputy voted Yes, No, or abstained for each legislation, which we represent as  $\{1, -1, 0\}$ , respectively. Since the temporal ordering of the voting sessions is likely to

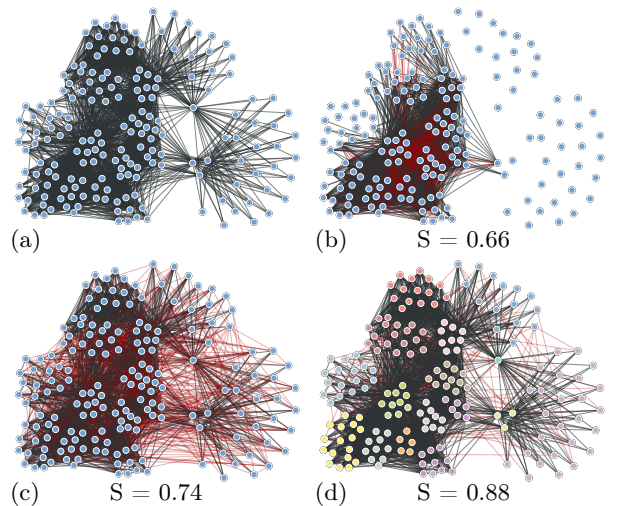
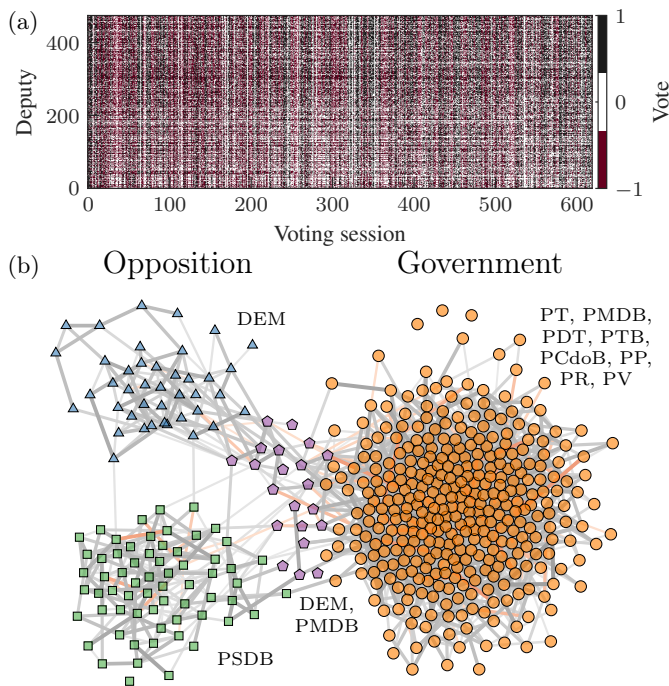


Figure 3. Reconstruction of a food web network [45] from  $M = 10^4$  samples of an Ising model at critical temperature. Edges marked in red are erroneous in the reconstruction. (a) Original network, (b) Optimal correlation thresholding, (c) Inverse correlations, (d) Joint reconstruction with community detection. The legends show the values of the posterior similarity (Eq. 11).

be of secondary importance to the voting outcomes, we assume the votes are sampled from an Ising model (the addition of zero-valued spins changes Eq. 9 only slightly by replacing  $2 \cosh(x) \rightarrow 1 + 2 \cosh(x)$ ). Fig 4 shows the result of the reconstruction, where the division of the nodes uncovers a cohesive government and a split opposition, as well as a marginal center group, which correlates very well with the known party memberships and can be used to predict unseen voting behavior (see Appendix D). In Fig 5 we show the result of the reconstruction of the directed network of influence between  $N = 1\,833$  twitter users from 58 224 re-tweets [49] using a SI epidemic model (the act of “re-tweeting” is modelled as an infection event, using Eqs. 5 and 6 with  $\gamma = 0$ ) and the nested DC-SBM. The reconstruction uncovers isolated groups with varying propensities to re-tweet, as well as groups that tend to influence a large fraction of users. By inspecting the geo-location metadata on the users, we see that the inferred groups amount to a large extent do different countries, although clear sub-divisions indicate that this is not the only factor governing the influence among users (see Appendix D 2).

*Conclusion* — We have presented a scalable Bayesian method to reconstruct networks from functional observations that uses the SBM as a structured prior, and, hence, performs community detection together with reconstruction. The method is nonparametric, and, hence, requires no prior stipulation of aspects of the network and size of the model, such as number of groups. By leveraging inferred correlations between edges, the SBM includes an additional source of evidence, and, thereby, improves the reconstruction accuracy, which in turn also



increases the accuracy of the inferred communities. The overall approach is general, requiring only appropriate functional model specifications, and can be coupled with an open ended variety of such models, other than those considered here.

Figure 4. Reconstruction of the interactions between members of the lower house of the Brazilian congress from the voting patterns of the 2007-2011 session, according to the Ising model. The node colors indicate the inferred groups. The edge thickness shows the posterior probability for each edge, and the color the magnitude of the coupling  $J_{ij}$ . The labels show the most frequent party membership for each group.

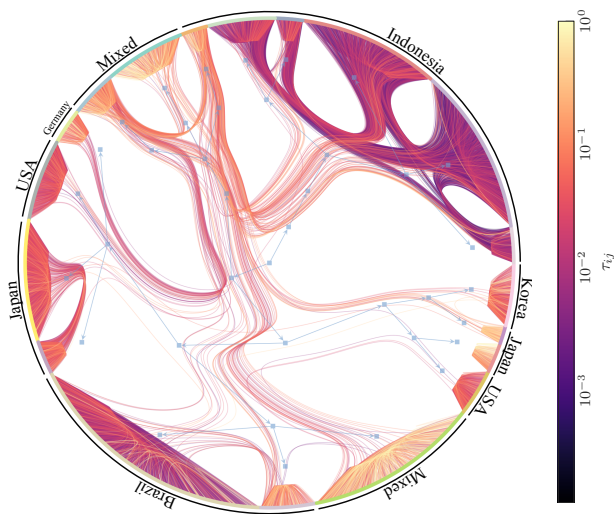


Figure 5. Reconstruction of the directed network of influence between  $N = 1833$  twitter users from 58 224 re-tweets, using a SI infection model. The hierarchical division represents the inferred fit of the nested DC-SBM (see Refs. [50, 51] for details on the layout algorithm), and the edge colors indicate the infection probabilities  $\tau_{ij}$  as shown in the legend. The text labels show the dominating country membership for the users in each group.



- 
- [1] Yong Wang, Trupti Joshi, Xiang-Sun Zhang, Dong Xu, and Luonan Chen, “Inferring gene regulatory networks from multiple microarray datasets,” *Bioinformatics* **22**, 2413–2420 (2006).
- [2] Michael Breakspear, “Dynamic models of large-scale brain activity,” *Nature Neuroscience* **20**, 340–352 (2017).
- [3] Matt J. Keeling and Pejman Rohani, “Estimating spatial coupling in epidemiological systems: a mechanistic approach,” *Ecology Letters* **5**, 20–29 (2002).
- [4] Nicolò Musmeci, Stefano Battiston, Guido Caldarelli, Michelangelo Puliga, and Andrea Gabrielli, “Bootstrapping Topological Properties and Systemic Risk of Complex Networks Using the Fitness Model,” *Journal of Statistical Physics* **151**, 720–734 (2013).
- [5] Eytan Bakshy, Itamar Rosenm, Cameron Marlow, and Lada Adamic, “The Role of Social Networks in Information Diffusion,” in *Proceedings of the 21st International Conference on World Wide Web, WWW ’12* (ACM, New York, NY, USA, 2012) pp. 519–528, event-place: Lyon, France.
- [6] Mark A. Kramer, Uri T. Eden, Sydney S. Cash, and Eric D. Kolaczyk, “Network inference with confidence from multivariate time series,” *Physical Review E* **79**, 061916 (2009).
- [7] Marc Timme, “Revealing Network Connectivity from Response Dynamics,” *Physical Review Letters* **98**, 224101 (2007).
- [8] Srinivas Gorur Shandilya and Marc Timme, “Inferring network topology from complex dynamics,” *New Journal of Physics* **13**, 013004 (2011).
- [9] Mor Nitzan, Jose Casadiego, and Marc Timme, “Revealing physical interaction networks from statistics of collective dynamics,” *Science Advances* **3**, e1600396 (2017).
- [10] Pieter Abbeel, Daphne Koller, and Andrew Y. Ng, “Learning Factor Graphs in Polynomial Time and Sample Complexity,” *Journal of Machine Learning Research* **7**, 1743–1788 (2006).
- [11] Guy Bresler, Elchanan Mossel, and Allan Sly, “Reconstruction of Markov Random Fields from Samples: Some Observations and Algorithms,” in *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques*, Lecture Notes in Computer Science (Springer, Berlin, Heidelberg, 2008) pp. 343–356.
- [12] Andrea Montanari and Jose A. Pereira, “Which graphical models are difficult to learn?” in *Advances in Neural Information Processing Systems 22*, edited by Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta (Curran Associates, Inc., 2009) pp. 1303–1311.
- [13] Holger Höfling and Robert Tibshirani, “Estimation of sparse binary pairwise markov networks using pseudolikelihoods,” *Journal of Machine Learning Research* **10**, 883–906 (2009).
- [14] H. Chau Nguyen, Riccardo Zecchina, and Johannes Berg, “Inverse statistical problems: from the inverse Ising problem to data science,” *Advances in Physics* **66**, 197–261 (2017).
- [15] Manuel Gomez Rodriguez, Jure Leskovec, and Andreas Krause, “Inferring Networks of Diffusion and Influence,” in *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’10* (ACM, New York, NY, USA, 2010) pp. 1019–1028.
- [16] Seth Myers and Jure Leskovec, “On the Convexity of Latent Social Network Inference,” in *Advances in Neural Information Processing Systems 23*, edited by J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta (Curran Associates, Inc., 2010) pp. 1741–1749.
- [17] Praneeth Netrapalli and Sujay Sanghavi, “Learning the Graph of Epidemic Cascades,” in *Proceedings of the 12th ACM SIGMETRICS/PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS ’12* (ACM, New York, NY, USA, 2012) pp. 211–222.
- [18] Chuang Ma, Han-Shuang Chen, Ying-Cheng Lai, and Hai-Feng Zhang, “Statistical inference approach to structural reconstruction of complex networks from binary time series,” *Physical Review E* **97**, 022301 (2018).
- [19] Bastian Prasse and Piet Van Mieghem, “Maximum-Likelihood Network Reconstruction for SIS Processes is NP-Hard,” arXiv:1807.08630 [physics] (2018), arXiv:1807.08630.
- [20] Braunstein Alfredo, Ingrosso Alessandro, and Muntoni Anna Paola, “Network reconstruction from infection cascades,” *Journal of The Royal Society Interface* **16**, 20180844 (2019).
- [21] Jakob Runge, Jobst Heitzig, Vladimir Petoukhov, and Jürgen Kurths, “Escaping the Curse of Dimensionality in Estimating Multivariate Transfer Entropy,” *Physical Review Letters* **108**, 258701 (2012).
- [22] Jie Sun, Dane Taylor, and Erik M. Bollt, “Causal Network Inference by Optimal Causation Entropy,” *SIAM Journal on Applied Dynamical Systems* (2015), 10.1137/140956166.
- [23] Zhesi Shen, Wen-Xu Wang, Ying Fan, Zengru Di, and Ying-Cheng Lai, “Reconstructing propagation networks with natural diversity and identifying hidden sources,” *Nature Communications* **5**, 4323 (2014).
- [24] Long Ma, Xiao Han, Zhesi Shen, Wen-Xu Wang, and Zengru Di, “Efficient Reconstruction of Heterogeneous Networks from Time Series via Compressed Sensing,” *PLOS ONE* **10**, e0142837 (2015).
- [25] Xiao Han, Zhesi Shen, Wen-Xu Wang, and Zengru Di, “Robust Reconstruction of Complex Networks from Sparse Data,” *Physical Review Letters* **114**, 028701 (2015).
- [26] Jingwen Li, Zhesi Shen, Wen-Xu Wang, Celso Grebogi, and Ying-Cheng Lai, “Universal data-based method for reconstructing complex networks with binary-state dynamics,” *Physical Review E* **95**, 032303 (2017).
- [27] Emily S. C. Ching, Pik-Yin Lai, and C. Y. Leung, “Reconstructing weighted networks from dynamics,” *Physical Review E* **91**, 030801 (2015).
- [28] Pik-Yin Lai, “Reconstructing network topology and coupling strengths in directed networks of discrete-time dynamics,” *Physical Review E* **95**, 022311 (2017).
- [29] Santo Fortunato and Darko Hric, “Community detection in networks: A user guide,” *Physics Reports* (2016), 10.1016/j.physrep.2016.09.002.
- [30] Tiago P. Peixoto, “Bayesian stochastic blockmodeling,” arXiv:1705.10225 [cond-mat, physics:physics, stat]

- (2017), arXiv: 1705.10225.
- [31] Quentin Berthet, Philippe Rigollet, and Piyush Srivastava, “Exact recovery in the Ising block-model,” arXiv:1612.03880 [math, stat] (2016), arXiv: 1612.03880.
- [32] Till Hoffmann, Leto Peel, Renaud Lambiotte, and Nick S. Jones, “Community detection in networks with unobserved edges,” arXiv:1808.06079 [physics] (2018), arXiv: 1808.06079.
- [33] M. E. J. Newman, “Network structure from rich but noisy data,” *Nature Physics* **14**, 542–545 (2018).
- [34] Tiago P. Peixoto, “Reconstructing Networks with Unknown and Heterogeneous Errors,” *Physical Review X* **8**, 041011 (2018).
- [35] Brian Karrer and M. E. J. Newman, “Stochastic block-models and community structure in networks,” *Physical Review E* **83**, 016107 (2011).
- [36] Tiago P. Peixoto, “Nonparametric Bayesian inference of the microcanonical stochastic block model,” *Physical Review E* **95**, 012317 (2017).
- [37] Tiago P. Peixoto, “Efficient Monte Carlo and greedy heuristic for the inference of stochastic block models,” *Physical Review E* **89**, 012804 (2014).
- [38] Roger Guimerà and Marta Sales-Pardo, “Missing and spurious interactions and the reconstruction of complex networks,” *Proceedings of the National Academy of Sciences* **106**, 22073–22078 (2009).
- [39] Romualdo Pastor-Satorras, Claudio Castellano, Piet Van Mieghem, and Alessandro Vespignani, “Epidemic processes in complex networks,” *Reviews of Modern Physics* **87**, 925–979 (2015).
- [40] Julian Besag, “Spatial Interaction and the Statistical Analysis of Lattice Systems,” *Journal of the Royal Statistical Society: Series B (Methodological)* **36**, 192–225 (1974).
- [41] Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, Augusta H. Teller, and Edward Teller, “Equation of State Calculations by Fast Computing Machines,” *The Journal of Chemical Physics* **21**, 1087 (1953).
- [42] W. K. Hastings, “Monte Carlo sampling methods using Markov chains and their applications,” *Biometrika* **57**, 97–109 (1970).
- [43] Aurelien Decelle, Florent Krzakala, Cristopher Moore, and Lenka Zdeborová, “Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications,” *Physical Review E* **84**, 066106 (2011).
- [44] Retrieved from [openflights.org](https://openflights.org).
- [45] Neo D. Martinez, “Artifacts or Attributes? Effects of Resolution on the Little Rock Lake Food Web,” *Ecological Monographs* **61**, 367–392 (1991).
- [46] Note that in this case our method also exploits the heterogeneous degrees in the network via the DC-SBM, which can in principle also aid the reconstruction, in addition to the community structure itself.
- [47] Refinements of this approach including TAP and BP corrections [14] yield the same performance for this example.
- [48] Aurélien Decelle and Federico Ricci-Tersenghi, “Pseudolikelihood Decimation Algorithm Improving the Inference of the Interaction Network in a General Class of Ising Models,” *Physical Review Letters* **112**, 070603 (2014).
- [49] Nathan O. Hodas and Kristina Lerman, “The Simple Rules of Social Contagion,” *Scientific Reports* **4**, 4343 (2014).
- [50] Tiago P. Peixoto, “Hierarchical Block Structures and High-Resolution Model Selection in Large Networks,” *Physical Review X* **4**, 011047 (2014).
- [51] D. Holten, “Hierarchical Edge Bundles: Visualization of Adjacency Relations in Hierarchical Data,” *IEEE Transactions on Visualization and Computer Graphics* **12**, 741–748 (2006).
- [52] Tiago P. Peixoto, “The `graph-tool` python library,” figshare (2014), 10.6084/m9.figshare.1164194, available at <https://graph-tool.skewed.de>.

## Appendix A: Nonparametric DC-SBM model summary

The DC-SBM used in this work is the same derived in detail in Ref. [36]. We give a succinct summary in the following. The marginal likelihood of the DC-SBM can be written as

$$P(\mathbf{A}|\mathbf{b}) = \int P(\mathbf{A}|\boldsymbol{\lambda}, \boldsymbol{\kappa}, \mathbf{b})P(\boldsymbol{\kappa}|\mathbf{b})P(\boldsymbol{\lambda}|\mathbf{b}) d\boldsymbol{\kappa} d\boldsymbol{\lambda}, \quad (\text{A1})$$

$$= P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b})P(\mathbf{k}|\mathbf{e}, \mathbf{b})P(\mathbf{e}|\mathbf{b}), \quad (\text{A2})$$

where  $\mathbf{e} = \{e_{rs}\}$  is the matrix of edge counts between groups, and  $\mathbf{k}$  is the degree sequence of the network, and with

$$P(\mathbf{A}|\mathbf{k}, \mathbf{e}, \mathbf{b}) = \frac{\prod_{r<s} e_{rs}! \prod_r e_{rr}!! \prod_i k_i!}{\prod_{i<j} A_{ij}! \prod_i A_{ii}!! \prod_r e_r!!}, \quad (\text{A3})$$

$$P(\mathbf{k}|\mathbf{e}, \mathbf{b}) = \prod_r \binom{n_r}{e_r}^{-1}, \quad (\text{A4})$$

$$P(\mathbf{e}|\mathbf{b}) = \bar{\lambda}^E / (\bar{\lambda} + 1)^{E+B(B+1)/2}, \quad (\text{A5})$$

being the microcanonical likelihood and corresponding noninformative priors. We further increase the explanatory power of this model [36] by replacing the microcanonical prior for the degrees with

$$P(\mathbf{k}|\mathbf{e}, \mathbf{b}) = P(\mathbf{k}|\boldsymbol{\eta})P(\boldsymbol{\eta}|\mathbf{e}, \mathbf{b}) \quad (\text{A6})$$

where  $\boldsymbol{\eta} = \{\eta_k^r\}$  are the degree frequencies of each group, with  $\eta_k^r$  being the number of nodes with degree  $k$  that belong to group  $r$ , and

$$P(\mathbf{k}|\boldsymbol{\eta}) = \prod_r \frac{\prod_k \eta_k^r!}{n_r!} \quad (\text{A7})$$

is a uniform distribution of degree sequences constrained by the overall degree counts, and finally

$$P(\boldsymbol{\eta}|\mathbf{e}, \mathbf{b}) = \prod_r q(e_r, n_r)^{-1} \quad (\text{A8})$$

is the distribution of the overall degree counts. The quantity  $q(m, n)$  is the number of different degree counts with the sum of degrees being exactly  $m$  and that have at most  $n$  non-zero counts, given by

$$q(m, n) = q(m, n-1) + q(m-n, n). \quad (\text{A9})$$

For the node partition we use the prior,

$$P(\mathbf{b}) = P(\mathbf{b}|\mathbf{n})P(\mathbf{n}|B)P(B) = \frac{\prod_r n_r!}{N!} \binom{N-1}{B-1}^{-1} N^{-1}. \quad (\text{A10})$$

which is agnostic to group sizes.

Finally, the hierarchical degree-corrected SBM (HDC-SBM) is obtained by replacing the uniform prior for  $P(\mathbf{e}|\mathbf{b})$  by a nested sequence of SBMs, where the edge counts in level  $l$  are generated by a SBM at a level above,

$$P(\mathbf{e}_l|\mathbf{e}_{l+1}, \mathbf{b}_l) = \prod_{r < s} \left( \binom{n_r^l n_s^l}{e_{rs}^{l+1}} \right)^{-1} \prod_r \left( \binom{n_r^l (n_r^l + 1)/2}{e_{rr}^{l+1}/2} \right)^{-1}, \quad (\text{A11})$$

where  $\binom{n}{m} = \binom{n+m-1}{m}$  is the multiset coefficient. The prior for the hierarchical partition is obtained using Eq. A10 at every level. The entire model above is also easily modified for directed networks. We refer to Ref. [36] for further details.

### Appendix B: Adapting multigraph models to simple graphs

The DC-SBM variations considered above generate multigraphs with self-loops, however the functional models presented in the main text operate on simple graphs. We amend this inconsistency in the same manner as in Ref. [34], by adapting the multigraph models to simple graphs in tractable way by generating multigraphs and then collapsing the multiple edges. In other words, if  $\mathbf{G}$  is a multigraph with entries  $G_{ij} \in \mathbb{N}$ , the collapsed simple graph  $\mathbf{A}(\mathbf{G})$  has binary entries

$$A_{ij}(G_{ij}) = \begin{cases} 1 & \text{if } G_{ij} > 0 \text{ and } i \neq j, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{B1})$$

Therefore, if  $\mathbf{G}$  is a multigraph generated by  $P(\mathbf{G}|\theta)$ , where  $\theta$  are arbitrary parameters, then the corresponding collapsed simple graph  $\mathbf{A}$  is generated by

$$P(\mathbf{A}|\theta) = \sum_{\mathbf{G}} P(\mathbf{A}, \mathbf{G}|\theta), \quad (\text{B2})$$

$$= \sum_{\mathbf{G}} P(\mathbf{A}|\mathbf{G})P(\mathbf{G}|\theta), \quad (\text{B3})$$

with

$$P(\mathbf{A}|\mathbf{G}) = \begin{cases} 1 & \text{if } \mathbf{A} = \mathbf{A}(\mathbf{G}), \\ 0 & \text{otherwise.} \end{cases} \quad (\text{B4})$$

Even if  $P(\mathbf{A}|\theta)$  cannot be computed in closed form, the joint distribution  $P(\mathbf{A}, \mathbf{G}|\theta) = P(\mathbf{A}|\mathbf{G})P(\mathbf{G}|\theta)$  is trivial, provided we have  $P(\mathbf{G}|\theta)$  in closed form. Therefore, instead of directly sampling from the posterior distribution

$$P(\mathbf{A}, \mathbf{b}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{A})P(\mathbf{A}, \mathbf{b})}{P(\mathcal{D})}, \quad (\text{B5})$$

we sample from the joint posterior

$$P(\mathbf{A}, \mathbf{G}, \mathbf{b}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{A})P(\mathbf{A}|\mathbf{G})P(\mathbf{G}, \mathbf{b})}{P(\mathcal{D})}, \quad (\text{B6})$$

using MCMC, treating the values  $G_{ij}$  as latent variables, and then we marginalize

$$P(\mathbf{A}, \mathbf{b}|\mathcal{D}) = \sum_{\mathbf{G}} P(\mathbf{A}, \mathbf{G}, \mathbf{b}|\mathcal{D}), \quad (\text{B7})$$

which is done simply by sampling from  $P(\mathbf{A}, \mathbf{G}, \mathbf{b}|\mathcal{D})$  and ignoring the actual magnitudes of the  $G_{ij}$  values, and the diagonal entries.

### Appendix C: Inference algorithm

The inference algorithm used here is identical to Ref. [34], with the only difference being the likelihoods for the forward model  $P(\mathcal{D}|\mathbf{A})$ . To summarize, we use MCMC to sample from the joint posterior distribution

$$P(\mathbf{A}, \mathbf{b}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{A})P(\mathbf{A}|\mathbf{b})P(\mathbf{b})}{P(\mathcal{D})}, \quad (\text{C1})$$

where  $\mathbf{b}$  is the partition of nodes used for the SBM. The MCMC algorithm consists of making proposals of the kind  $P(\mathbf{b}'|\mathbf{A}, \mathbf{b})$  and  $P(\mathbf{A}'|\mathbf{A}, \mathbf{b})$  for the partition and network, respectively (or equivalently for any other remaining model parameter), and accepting them according to the Metropolis-Hastings probability

$$\min \left( 1, \frac{P(\mathbf{A}', \mathbf{b}'|\mathcal{D})P(\mathbf{A}|\mathbf{A}', \mathbf{b}')P(\mathbf{b}|\mathbf{A}', \mathbf{b}')}{P(\mathbf{A}, \mathbf{b}|\mathcal{D})P(\mathbf{A}'|\mathbf{A}, \mathbf{b})P(\mathbf{b}'|\mathbf{A}, \mathbf{b})} \right), \quad (\text{C2})$$

which does not require the computation of the intractable normalization constant  $P(\mathcal{D})$ . In practice, at each step in the chain we make either a move proposal for  $\mathbf{A}$  or  $\mathbf{b}$ , not both at once. For the node partition, we use the move proposals described in Refs. [36, 37], where for any given node  $i$  in group  $r$  we propose to move it to group  $s$  (which can be previously unoccupied, in which case it is labelled  $s = B + 1$ ) according to

$$P(b_i = r \rightarrow s|\mathbf{A}, \mathbf{b}) = d\delta_{s, B+1} + (1-d)(1-\delta_{s, B+1}) \sum_{t=1}^B P(t|i) \frac{e_{ts} + \epsilon}{e_t + \epsilon B}, \quad (\text{C3})$$

where  $P(t|i) = \sum_j A_{ij} \delta_{b_j, t} / k_i$  is the fraction of neighbours of  $i$  that belong to group  $t$ ,  $\epsilon > 0$  is a small parameter which guarantees ergodicity, and  $d$  is the probability of moving to a previously unoccupied group. (If  $k_i = 0$ , we assume  $P(b_i = r \rightarrow s|\mathbf{A}, \mathbf{b}) = d\delta_{s, B+1} + (1-d)(1 - \delta_{s, B+1})/B$ .) This move proposal attempts to use the currently known large-scale structure of the network to better inform the possible moves of the node, without biasing with respect to group assortativity. The parameters  $d$  and  $\epsilon$  do not affect the correctness of the algorithm,



only the mixing time, which is typically not very sensitive, provided they are chosen within a reasonable range (we used  $d = 0.01$  and  $\epsilon = 1$  throughout). When using the HDC-SBM, we used the variation of the above for hierarchical partitions described in Ref. [36]. The move proposals above require only minimal bookkeeping of the number edges incident on each group, and can be made in time  $O(k_i)$ , which is also the time required to compute the ratio in Eq. C2, independent on how many groups are currently occupied.

For the network, we change the values of the latent multigraph  $\mathbf{G}$  with unit proposals

$$P(G'_{ij} = G_{ij} + \delta | \mathbf{G}) = \begin{cases} 1/2 & \text{if } G_{ij} > 0, \\ 1 & \text{if } G_{ij} = 0 \text{ and } \delta = 1, \\ 0 & \text{otherwise,} \end{cases} \quad (\text{C4})$$

for  $\delta \in \{-1, 1\}$ . We choose the entries to update with a probability given by the current DC-SBM,

$$P(i, j | \mathbf{G}, \mathbf{b}) = \kappa_i \kappa_j m_{b_i, b_j}, \quad (\text{C5})$$

with

$$\kappa_i = \frac{k_i + 1}{\sum_j \delta_{b_j, b_i} k_j + 1} \quad (\text{C6})$$

being the probability of selecting node  $i$  from its group  $b_j$ , proportional to its current degree plus one, and

$$m_{rs} = \frac{e_{rs} + 1}{\sum_{tu} e_{rs} + 1} \quad (\text{C7})$$

is the probability of selecting groups  $(r, s)$ , where  $e_{rs} = \sum_{ij} G_{ij} \delta_{b_i, r} \delta_{b_j, s}$ . The above probabilities guarantee that every entry will be eventually sampled, but it tends to probe denser regions more frequently, which we found to typically lead to faster mixing times. This sampling can be done in time  $O(1)$ , simply by keeping urns of vertices and edges according to the group memberships. The time required to compute the ratio in Eq. C2 is also  $O(1)$  for the DC-SBM and  $O(L)$  for the HDC-SBM, where  $L$  is the hierarchy depth, again independent of the number of occupied groups.

## 1. Algorithmic complexity

When combining the move proposals defined above for the partition and network, the time required to perform  $N$  node move proposals and  $E$  edge addition or removal proposals is  $O(\langle k \rangle N + EM)$ , where  $\langle k \rangle$  is the average degree, and  $M$  is the number of samples per node of the functional model (i.e. SIS or Ising). The  $O(EM)$  contribution is seen by noting that the addition and removal of an edge requires the re-computation of the likelihood  $P(\mathcal{D} | \mathbf{A})$  involving only terms associated with each endpoint over all  $M$  samples, each requiring only  $O(1)$  computations. For the SIS model we note that we need only

to keep track of the summary quantities  $m_i(t)$  for each node, and update them by adding or subtracting contributions for each added or removed edge, and the same is true for the Ising model with respect to edge contributions to the Hamiltonian. This linear complexity of sweeps allows for the reconstruction of large networks.

For dynamical data where changes of the state of each node are relatively rare (e.g. in a SI or SIR dynamics, a node changes its state only once or twice, respectively, for the whole cascade), it is possible to optimize the inference algorithm by listing for each node only its initial state and the times it changes, together with the new state values. In this way, the contribution to the likelihood of a single node can be computed by only going through the times that its neighbours or the node itself change state, instead of the whole time-span of the dynamics. Therefore, the complexity for a whole MCMC sweep changes to  $O(\langle k \rangle N + Ea)$ , where  $a$  is the average number of times a single node changes its state during the whole dynamics. For very active dynamics we have  $a = O(M)$ , and hence this algorithm has the same complexity as the version above, but for  $a \ll O(M)$  it gives noticeable speed-ups.

In addition to the algorithmic complexity of each sweep, the MCMC needs time to converge to the target posterior distribution. This mixing time depends not only on the structure of the network being reconstructed, how easy it is to uncover it from the data, but also on how close the algorithm is initiated to the target distribution. Because of this, it is not straightforward to estimate the general algorithmic complexity of the mixing time. In our numerical experiments we found that both starting from random or empty networks lead to reasonable mixing times in most cases, and the results coincide with initializing from the true planted network (which as expected, shows faster equilibration).

A reference implementation of the above algorithm is freely available as part of the `graph-tool` library [52].

## Appendix D: Datasets with empirical dynamics

### 1. Cross-validation

To further evaluate the ability of our reconstruction method to capture the empirical voting behavior of deputies in the lower house of the Brazilian congress, we randomly divided all  $M = 619$  voting sessions into  $M - M_t$  “training” sessions, used to fit the model, and  $M_t$  test sessions, used to compare with the predictions from the model. To evaluate the prediction error, the correlation matrix  $C_{ij} = \langle s_i s_j \rangle - \langle s_i \rangle \langle s_j \rangle$  was computed

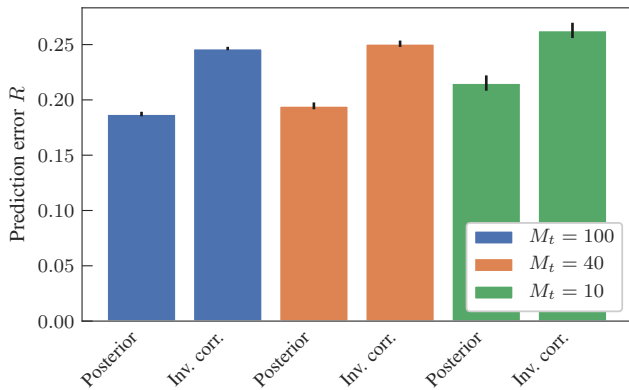


Figure 6. Cross-validation results for the voting behavior of deputies in the lower house of the Brazilian congress, comparing the posterior sampling approach considered in the main text and the inverse correlation method.

Group	User countries
1	Brazil $\times$ 216, Japan $\times$ 2, Germany $\times$ 2, USA $\times$ 1, Russia $\times$ 1, Argentina $\times$ 1, Colombia $\times$ 1, France $\times$ 1, South Africa $\times$ 1
2	Japan $\times$ 104, UK $\times$ 1, China $\times$ 1, Italy $\times$ 1, USA $\times$ 1
3	Indonesia $\times$ 42, India $\times$ 10, UK $\times$ 8, Germany $\times$ 6, USA $\times$ 6, Thailand $\times$ 1, Brazil $\times$ 1, Bali $\times$ 1, Argentina $\times$ 1, Australia $\times$ 1
4	Indonesia $\times$ 52, Germany $\times$ 6, UK $\times$ 4, USA $\times$ 3, Russia $\times$ 2, Australia $\times$ 2, New Zealand $\times$ 1, India $\times$ 1, Botswana $\times$ 1, Brazil $\times$ 1
5	USA $\times$ 54, UK $\times$ 4, Russia $\times$ 1, Japan $\times$ 1, Indonesia $\times$ 1
6	USA $\times$ 23, Brazil $\times$ 9, UK $\times$ 6, Netherlands $\times$ 3, Russia $\times$ 2, Canada $\times$ 2, Italy $\times$ 2, Germany $\times$ 2, Mexico $\times$ 1, South Africa $\times$ 1, Kuwait $\times$ 1, Austria $\times$ 1, Romania $\times$ 1, Finland $\times$ 1, Japan $\times$ 1, Philippines $\times$ 1, Egypt $\times$ 1, Argentina $\times$ 1, Chile $\times$ 1
7	USA $\times$ 18, India $\times$ 8, Portugal $\times$ 2, Brazil $\times$ 2, Japan $\times$ 2, UK $\times$ 2, Guernsey $\times$ 1, Malaysia $\times$ 1, Mexico $\times$ 1, Australia $\times$ 1, Russia $\times$ 1, France $\times$ 1, Vietnam $\times$ 1, Spain $\times$ 1, Venezuela $\times$ 1, Philippines $\times$ 1
8	Japan $\times$ 35
9	Brazil $\times$ 31, USA $\times$ 2, UK $\times$ 1
10	Korea $\times$ 24, USA $\times$ 1, Argentina $\times$ 1, Russia $\times$ 1, Japan $\times$ 1
11	Indonesia $\times$ 22, UK $\times$ 2, Australia $\times$ 1, India $\times$ 1, USA $\times$ 1
12	USA $\times$ 11, Japan $\times$ 2, Germany $\times$ 2, France $\times$ 1, Korea $\times$ 1, Thailand $\times$ 1, Brazil $\times$ 1, Chile $\times$ 1, Indonesia $\times$ 1
13	USA $\times$ 4, Indonesia $\times$ 4, India $\times$ 4, UK $\times$ 3, Australia $\times$ 2, Canada $\times$ 1
14	Japan $\times$ 8, Venezuela $\times$ 6, USA $\times$ 1, Chile $\times$ 1
15	Indonesia $\times$ 14, Japan $\times$ 1
16	Indonesia $\times$ 7, USA $\times$ 3, Turkey $\times$ 1, Philippines $\times$ 1, Brazil $\times$ 1
17	USA $\times$ 9, Canada $\times$ 1, France $\times$ 1, Belgium $\times$ 1
18	Germany $\times$ 5, USA $\times$ 1, Russia $\times$ 1, France $\times$ 1
19	Japan $\times$ 4

Table I. Country memberships of twitter users, according to the groups inferred by the reconstruction method.

Group	Parties
1	PMDB $\times$ 83, PT $\times$ 70, PP $\times$ 38, PR $\times$ 36, PSB $\times$ 27, PDT $\times$ 22, PTB $\times$ 17, PV $\times$ 14, PCdoB $\times$ 12, PSC $\times$ 10, DEM $\times$ 5, PMN $\times$ 5, PRB $\times$ 3, PSOL $\times$ 3, PHS $\times$ 2, PTdoB $\times$ 1, PTC $\times$ 1
2	PSDB $\times$ 54, PPS $\times$ 8, PFL $\times$ 2
3	DEM $\times$ 41
4	DEM $\times$ 8, PMDB $\times$ 6, PPS $\times$ 4, PP $\times$ 2, PSB $\times$ 1

Table II. Party affiliation of deputies of the lower house of the Brazilian congress, according to the groups they were classified by the reconstruction method. Parties in red belong to the center-left government coalition, and in blue to the right-wing opposition.

from the posterior distribution of the fitted model via

$$\langle s_i s_j \rangle = \sum_{\mathbf{A}, \mathbf{b}, \beta, \mathbf{J}, \mathbf{h}} s_i s_j P(\mathbf{s} | \mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) P(\mathbf{A}, \mathbf{b}, \beta, \mathbf{J}, \mathbf{h} | \bar{\mathbf{s}}_t) \quad (\text{D1})$$

$$\langle s_i \rangle = \sum_{\mathbf{A}, \mathbf{b}, \beta, \mathbf{J}, \mathbf{h}} s_i P(\mathbf{s} | \mathbf{A}, \beta, \mathbf{J}, \mathbf{h}) P(\mathbf{A}, \mathbf{b}, \beta, \mathbf{J}, \mathbf{h} | \bar{\mathbf{s}}_t) \quad (\text{D2})$$

where  $\bar{\mathbf{s}}_t$  is the training data, and compared with the correlation matrix obtained for the test data  $C_{ij}^{(t)} = \langle s_i s_j \rangle_t - \langle s_i \rangle_t \langle s_j \rangle_t$ , via

$$\langle s_i s_j \rangle_t = \frac{1}{M_t} \sum_m s_i^m s_j^m, \quad \langle s_i \rangle_t = \frac{1}{M_t} \sum_m s_i^m, \quad (\text{D3})$$

where the sums go over microstates of the test data. The prediction error is then computed as

$$R = \frac{1}{\binom{N}{2}} \sum_{i < j} |C_{ij} - C_{ij}^{(t)}|. \quad (\text{D4})$$

We repeated the calculation using  $M_t \in \{10, 40, 100\}$ , and for each value of  $M_t$  we averaged the results over 30 random choices of the test data. We also compared with the reconstruction obtained via inverse correlations [14]. The results are shown in Fig. 6. As can be seen, the Bayesian joint reconstruction method outperforms the results based on inverse correlations. The prediction error decreases with larger  $M_t$ , as in this limit the fluctuations in the dynamics are averaged out.

## 2. Comparison with metadata

Here we expand on the comparison of the community structure found both for the Brazilian congress as well as the twitter data, with metadata available in both cases.

In table II we list the party affiliations of the deputies according to the groups they were classified by the method. The largest group accounts for all left-wing parties as well as the center parties belonging to the government coalition, whereas groups 2 and 3 accounts for the

right-wing opposition. Group 4 is composed of a small number deputies who are members of both government and opposition parties, but vote independently.

In table I is shown the country of each twitter user, independently obtained via twitter's API (not contained in the original dataset of Ref. [49]), according to each group identified by the reconstruction method. As can

be seen, most groups are characterized by a single dominating country, with only a few exceptions. This indicates, plausibly, that the probability of re-tweets is largely shaped by language and cultural barriers. Nevertheless, the method also uncovers subdivisions within the distinct geographical locations, indicating that this is not the only factor determining the influence among users.