# Learning of networked spreading models from noisy and incomplete data

Mateusz Wilinski and Andrey Y. Lokhov

*Theoretical Division, Los Alamos National Laboratory, Los Alamos, USA*

Recent years have seen a lot of progress in algorithms for learning parameters of spreading dynamics from both full and partial data. Some of the remaining challenges include model selection under the scenarios of unknown network structure, noisy data, missing observations in time, as well as an efficient incorporation of prior information to minimize the number of samples required for an accurate learning. Here, we introduce a universal learning method based on scalable dynamic message-passing technique that addresses these challenges often encountered in real data. The algorithm leverages available prior knowledge on the model and on the data, and reconstructs both network structure and parameters of a spreading model. We show that a linear computational complexity of the method with the key model parameters makes the algorithm scalable to large network instances.

## I. INTRODUCTION

Spreading models are routinely used to generate predictions for a plethora of diffusion processes on networks, whereby infectious diseases, opinions, or failures propagate in natural, social, and technological systems [1, 2]. In these models, the nodes typically go from inactive to an active state through interactions with their neighbors on a network, similarly to how an infection is passed from one person to another. When the model structure or parameters are unknown, it is natural to consider the inverse problem of learning of the spreading model from data. Available data usually takes form of reported activation times for nodes in a network in several observed activation cascades. Whether due to a limited observation budget or imperfect reporting, the accessible data is unlikely to be perfect, and may be subject to uncertainty or provide only partial information on the system. This motivates the design of robust methods for selecting a spreading model from incomplete and noisy data. In this paper, we introduce a scalable learning method that addresses the challenges related to imperfect data.

Reconstruction of networked spreading models has been addressed by a number of works in recent years. In the context of full and exact observation of nodes' activation times, [3–6] showed that maximum-likelihood type approaches succeed to learn the model structure and parameters. The problem becomes significantly more difficult when the system is only partially observed, i.e., only a fraction of nodes report information about their activation times. In this case, the maximum likelihood approach has exponential complexity with respect to the number of unobserved nodes, which warrants learning methods based on the direct problem of predicting the model dynamics [7]. This approach has lead to an efficient algorithm for learning the parameters of the spreading model on a known network, based on minimization of the distance between observed and model-predicted node marginal probabilities [8]. This choice is motivated by the fact that marginal probabilities of nodes' activation can be estimated in a computationally efficient way using a method known as dynamic message passing (DMP) [9–28]. DMP is an inference method for spreading processes on networks derived from a classical belief propagation (BP) algorithm [29, 30], and thus showing similar properties to BP [17, 21]. The key result of [8] stated that the model learned using the minimization of distance between marginal probabilities generates better predictions compared to the model with ground-truth parameters when DMP is used as the inference algorithm. This has been shown through an empirical study on a popular spreading model known as Independent Cascade [31, 32], equivalent to a popular epidemic spreading susceptible-infected-recovered (SIR) model [33, 34] with deterministic recovery. Due to a prediction-centric focus, the algorithm proposed in [8] has been referred to as SLICER (Scalable Learning of Independent Cascade Effective Representation).

On many sparse network instances, SLICER produces high-quality estimates of model parameters in a time scaling linearly with the system size, even for fraction of hidden nodes up to 25% of the network. However, from the results of [8], it is not immediately obvious what the expected accuracy limits of the method are, i.e., what error can be expected when the number of unobserved nodes reaches very high values. In order to get insight into the performance of SLICER, in section III A we run a systematic evaluation of the performance of SLICER for learning of the spreading model parameters on a variety of random graph classes under diverse fractions of unobserved nodes.

Partial node-observability of the system is only one of the several ways in which the data may be incomplete or uncertain. Some of the challenges often encountered in real data include model selection under the scenarios of unknown or partially known network structure, noisy data, missing observations in time, or necessity to include prior information in order to minimize the required amount of data. Scenarios that we consider are presented in Figure 1. The first challenge that we address is the lack of information about the network structure. This task has been recently studied in [35], however, the proposed method involves cubic complexity in the number of nodes, making its application to large networks prohibitive. The second challenge deals with incorpo-
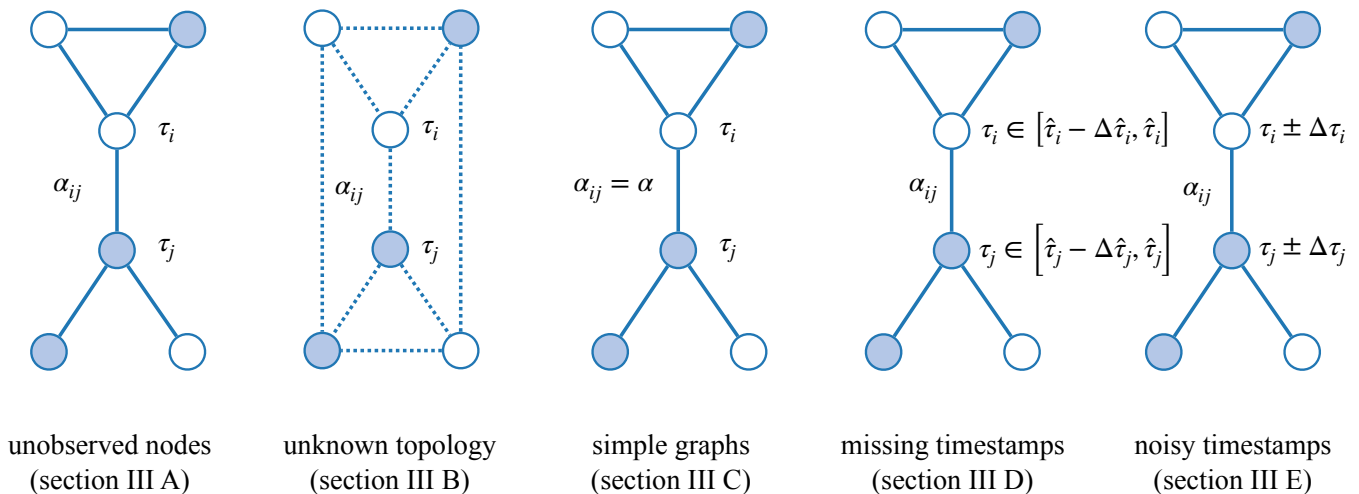
FIG. 1: Schematic representation of different scenarios of incomplete and uncertain data that we consider in learning of networked spreading models. We focus on the general setting where a fraction of the nodes never reports information (depicted as empty nodes). In addition, we treat cases where network structure is not known or known only partially; where the number of data samples is small, but prior information on the parameters is available; and where the observed node activation timestamps are missing or noisy. At the end of our experimental evaluation, we address learning of real-world network instances under a combination of these scenarios.

ration of prior information on the model parameters in the case where the number of available samples is small. For instance, in epidemiological applications, at most a few realizations of the dynamics can be observed, but the spreading is often modeled by a single parameter (transmission probability) [36]. We show how to leverage this prior information and learn the simple graph models from a few observed trajectories. The third challenge is missing data in time, where nodes report their activation times only during limited time windows. This task was analysed in [37], but with an additional assumption that the full probabilistic trace for each node is available. A likelihood-free approach for the missing time case was proposed in [38], together with additional goal of finding the source of spreading. A variant of the scenario of missing timestamps assumes that only the complete final state of the dynamics is available. One of the first results for this variant of the problem assumed that the unknown network is a tree [39], but later this assumption has been relaxed in [40–42]. In our work, we show that the spreading model can be accuractly recovered with minimal available information in the temporal space, even in the case of partial node-observability. The last challenge that we consider, often related to the reporting procedure, is uncertainty in timestamps, which may be modeled as an additional noise added to the observed data. Similar setting was analysed in [43, 44], but not paired with any other type of missing information such as partial node-wise reporting.

As we show in this work, in most of these scenarios, the objective function used in SLICER method is not directly applicable, resulting in biased estimates of the model pa-

rameters. Our overreaching goal consists in generalizing SLICER, introducing a universal learning framework capable of addressing all these challenges. Here, we show how to crucially define the objective function in such a way that it incorporates the available prior information, and provides a high-quality reconstruction of the model parameters. As demonstrated below and contrary to the solution of these challenges on a case-by-case basis, this universality presents an advantage in that data uncertainties present in the data and pertaining to different challenges can be treated at the same time. As a demonstration of such a universality, in the end of our study, we illustrate the approach on real-world network instances under a combination of these challenges.

## II. METHODS

In the section, we first pose the learning problem using the Independent Cascade (IC) model as the dynamics of choice. Then, we discuss the dynamic message passing (DMP) inference method that lies at the foundation of our approach. Finally, we explain the details of our learning method.

### A. Problem formulation: Learning of the Independent Cascade model

In this paper, we focus on the Independent Cascade model [31, 32]. Define the spreading network as a graph $G = (V, E)$, consisting of the set of vertices $V = \{v_i\}$,

and the set of edges $E = \{(v_i, v_j) \mid v_i, v_j \in V\}$. Each node $v_i$ in the IC model can be in one of two states: active and inactive. Assuming node-wise independent initial condition for the states of the nodes at time $t = 0$, which can be either deterministic or stochastic, further dynamics is subject to a single stochastic rule. If node $i$ gets activated at any time $t$, it has only one chance to activate any of its inactive neighbors $j$ at time $t+1$. This happens independently for each of the inactive neighbor with edge dependent probability $\alpha_{ij}$:

$$A(i) + I(j) \xrightarrow{\alpha_{ij}} A(i) + A(j). \quad (1)$$

Regardless of activating any of its neighbors, node $i$ remains active forever, but it cannot activate any other node in the future beyond $t + 1$. In what follows, we assume that $G$ is undirected and there are no self-loops or multi-links (although this assumption can be relaxed in general). The dynamics continues until a predefined number of steps $T$. A single realisation of this process is called a cascade. Since each node can be activated only once, every cascade $c$ is fully described by a set of activation times $\{\tau_i^c\}_{i \in V}$. Additionally, if node $i$ does not get activated before or at the specified time $T$ in cascade $c$, we assign $\tau_i^c = *$ as its activation time. It means that $*$ summarizes all future activity of a given node. More details about this notation is given in the Appendix A. In all our experiments, we use a single seed as an initial condition (this setting corresponds to the most popular scenario, but can be relaxed as well). This means that there is only one node activated at the beginning of the spreading process, although it may be a different node for each specific cascade. The number of observed cascades is denoted by $M$.

The problem that we consider in this work is as follows: given a set of observed activation times $\{\tau_i^c\}_{i \in O, c \in C}$ where $O \subset V$ is the set of observed nodes and $C$ is the set of cascades, as well as using additional information on the cascades (e.g., prior knowledge that $\tau_i^c$ are noisy), learn model parameters, i.e., estimate $\{\alpha_{ij}\}_{(ij) \in V \times V}$ so that they are close to the parameters of the ground-truth model, denoted as $\{\alpha_{ij}^*\}_{(ij) \in E}$. In the case of the unknown network, thresholded values of $\{\alpha_{ij}\}_{(ij) \in V \times V}$ away from small values close to zero define the recovered network structure $\widehat{E}$ that can be compared to the ground-truth set of edges $E$. In what follows, the fraction of hidden nodes $|V \backslash O|/|O|$ is denoted as $\xi$.

We chose the IC model for simplicity reasons. On one hand, it does capture basic properties of many known spreading processes, including a possibility of an early cascade termination. On the other, it simplifies the analytical equations presented in the next subsection, making the proposed approach easier to follow and understand. Our approach can be can easily be generalized to a broad class of more complex spreading models on networks for which DMP equations (explained next) are known.

## B. Inference method: Dynamic Message Passing

As discussed in the Introduction, presence of partial information warrants learning methods that feature inference algorithms for predicting the model dynamics as a subroutine. One of the key observables that quantify the spread is given by the *influence function*: the number of expected number of activated nodes at a certain time $t$ for a given initial condition. For a fixed initial conditions in cascade $c$, the influence function is given by the sum of marginal probabilities $p_i^c(t)$ of activation of each node $i \in V$ at at time $t$ [20]. Prediction of influence function or marginal probabilities is known to be #P-hard [45, 46], and hence approximate methods need to be used. A classical way of estimating the influence function consists in using Monte-Carlo simulations. However, this approach typically requires a large sampling factor to provide a reliable estimate [32, 47–51]. Aiming at an accelerated approach that saves this sampling factor, we use the inference method known as Dynamic Message Passing [9–28]. For IC model [21], DMP estimates marginal probabilities of activation in a linear time in both system size and duration of the dynamics, and has the properties of being exact on graphs without loops and asymptotically exact on random graphs, providing an upper-bound of the influence function on general networks.

For the IC model, the equations take the following form:

$$p_i^c(t) = 1 - \left(1 - \bar{p}_i^c\right) \prod_{k \in \partial i} \left(1 - \alpha_{ki} \cdot p_{k \to i}^c(t-1)\right), \quad (2)$$

$$p_{j \to i}^c(t) = 1 - \left(1 - \bar{p}_j^c\right) \prod_{k \in \partial j \backslash i} \left(1 - \alpha_{kj} \cdot p_{k \to j}^c(t-1)\right), \quad (3)$$

where $p_i^c(t)$ is the marginal probability of node $i$ being active at time $t$ under initial conditions of cascade $c$, $p_{i \to j}^c(t)$ is the same probability, but on an auxiliary graph where node $j$ was removed, $\partial i$ denotes the set of neighbors of node $i$ in the graph $G$, and $\partial j \backslash i$ denotes the set of neighbors of node $j$ in the graph $G$ except $i$. We also denote an initial condition for node $i$ in cascade $c$ as $\bar{p}_i^c = p_i^c(0)$. Note that for known model parameters, the marginals given by DMP depend only on the initial condition $\bar{p}_i^c$. Detailed properties of the above equations were studied in [21]. In Appendix A, following the approach of [17], we provide an alternative derivation of DMP equations for IC model that connects them and their properties to the classical belief propagation algorithm [29, 30].

## C. Learning Framework

In this section, we explain the details of our method. We start by stating the approach introduced in [8] and known as Scalable Learning of Independent Cascade Effective Representation (SLICER), which will be used as

a baseline method in all numerical experiments below. Subsequently, we will show that the objective used in SLICER need to be modified to account for the prior information on the uncertainty in the data, and enhanced with the known constraints on the model parameters – an approach referred to as SLICER+.

Let us focus on the case of perfectly observed information from a subset of visible nodes $O$ [8]. SLICER is based on the minimization of the Kullback–Leibler divergence between the empirical marginal distributions computed from the data, and the marginals estimated from the model using DMP, on the observed nodes. Marginal distributions depend on non-local model parameters, including those adjacent to hidden nodes, and thus allowing for the reconstruction of the entirety of model parameters under sufficient observations. Minimization of KL distance on marginals on the observed nodes in our case is equivalent to maximizing:

$$\mathcal{O} = \sum_{c \in C} \sum_{i \in O} \log \mu_i^c(\tau_i^c), \qquad (4)$$

where $C$ is the set of available cascades, $O$ is the set of visible nodes and $\mu_i^c(t)$ is the marginal probability of node $i$ being activated under cascade $c$ precisely at time $t$. These marginal probabilities can be calculated based on the marginal variables $p_i^c(t)$ used in the DMP equations above:

$$\mu_i^c(t) = p_i^c(t) \cdot \mathbb{1}_{(t<T)} - p_i^c(t-1) \cdot \mathbb{1}_{(t>0)} + \mathbb{1}_{(t=T)}, \quad (5)$$

where $\mathbb{1}$ stands for the indicator function.

The cost function (4) has been first proposed in [7] showing an asymptotic consistency of DMP-based recovery using this objective, but providing an inefficient optimization algorithm. The work [8] proposed an efficient algorithm for minimizing (4), which gave rise to the SLICER algorithm.

For a given set of parameters $\alpha_{ij}$, marginal probabilities depend only on the initial condition. Therefore, we can re-write the objective in the following way:

$$\mathcal{O} = \sum_{s \in S} \sum_{i \in O} \sum_{\tau_i^s} m^{\tau_i^s} \log \mu_i^s(\tau_i^s), \qquad (6)$$

where $S$ is a set of all the initial conditions occurring across all the cascades $C$ and $m^{\tau_i^s}$ is the number of times node $i$ gets activated at time $\tau_i^s$ under initial condition $s$. This equivalent reformulation allows one to reduce the computing cost by evaluating DMP equations only $|S|$ times instead of $M = |C|$. In our simulations, where $s$ is assumed to be a single node, it means that DMP will not be run more than $N$ times, regardless of the number of available cascades. In the case of a stochastic initial condition of all of the cascades, e.g., when each node is independently initially activated with probability $\frac{1}{N}$, the size of $S$ is $|S| = 1$, making the whole computation far less costly.

In order to maximize (6) we use a Lagrangian formulation, where the objective function $\mathcal{O}$ (6) is supplemented

with constrains on marginal provided via DMP,

$$\mathcal{L} = \underbrace{\mathcal{O}}_{objective} + \underbrace{\mathcal{C}}_{constraints}, \qquad (7)$$

where in the absence of any additional prior information on the parameters, $\mathcal{C}$ are given by DMP equations (2)-(3) reweighted by Lagrange multipliers $\lambda_i^s(t)$ for all nodes $i \in V$ and $\lambda_{i \to j}^s(t)$ for all edges $(ij) \in E$ of graph $G$ for each time $t$ and for each cascade in the class $s$, see Appendix B for details. The iterative equations constituting SLICER is obtained by differentiating the Lagrangian (7) with respect to all variables: $p_i^s(t)_{i \in V}$, $p_{i \to j}^s(t)_{(ij) \in V \times V}$, $\alpha_{ij\,(ij) \in V \times V}$, $\lambda_i^s(t)_{i \in V}$, and $\lambda_{i \to j}^s(t)_{(ij) \in V \times V}$. This results in DMP equations (2)-(3) in the primal space; DMP-like equations on the Lagrange multipliers $\lambda_i^s(t)_{i \in V}$, and $\lambda_{i \to j}^s(t)_{(ij) \in V \times V}$ in the dual space; and update equations for the parameters $\alpha_{ij\,(ij) \in V \times V}$, see Appendix B.

When prior information about the parameters is available, it can be directly incorporated in the term $\mathcal{C}$. As an example, the constraint can be greatly simplified if we know that for each $(ij) \in E$, $\alpha_{ij} = \alpha$ (the case of so-called simple graphs, representing a popular case in epidemiological models). As we show below in the section III C, the gradient of the Lagrangian with respect to the model parameter $\alpha$ in this case reads for $\alpha \neq 0$:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = -\frac{1}{\alpha} \sum_{s \in S} \sum_{t=0}^{T-1} \sum_{(i,j) \in V \times V} \lambda_{i \to j}^s(t) \cdot p_{i \to j}^s(t). \qquad (8)$$

The gradient in Eq. (B6) can be used to learn parameter $\alpha$ using an iterative procedure with $\alpha \longleftarrow \alpha + \varepsilon \frac{\partial \mathcal{L}}{\partial \alpha}$, where $\varepsilon$ is a learning rate. Lagrangian formulation ensures that a single gradient descent step has the worst-case computational complexity $O(|E||T||S|)$, which means it is linear in the system size, cascade length and the number of initial conditions.

In a similar manner, under different scenarios of incomplete data considered in Figure 1, the loss function $\mathcal{O}$ can be appropriately modified to include the prior information $A$ on the type of the uncertainty in the data, as we discuss in the next section below. Collectively, we refer to the algorithm based on modified objective $\mathcal{O}$ and constraints $\mathcal{C}$ as to SLICER+. The full derivation of SLICER+ for each of the challenges in Figure 1 is given in the Appendix B. Efficient implementation of the algorithm, taking advantage of the DMP-like equations on the evolution of the Lagrange multiplies in the dual space and assuring the linear complexity of the algorithm, is available in the Appendix C.

## III. RESULTS

In this section we present details on the SLICER+ approach designed for dealing with different types of challenges regarding uncertainty and partial observability of
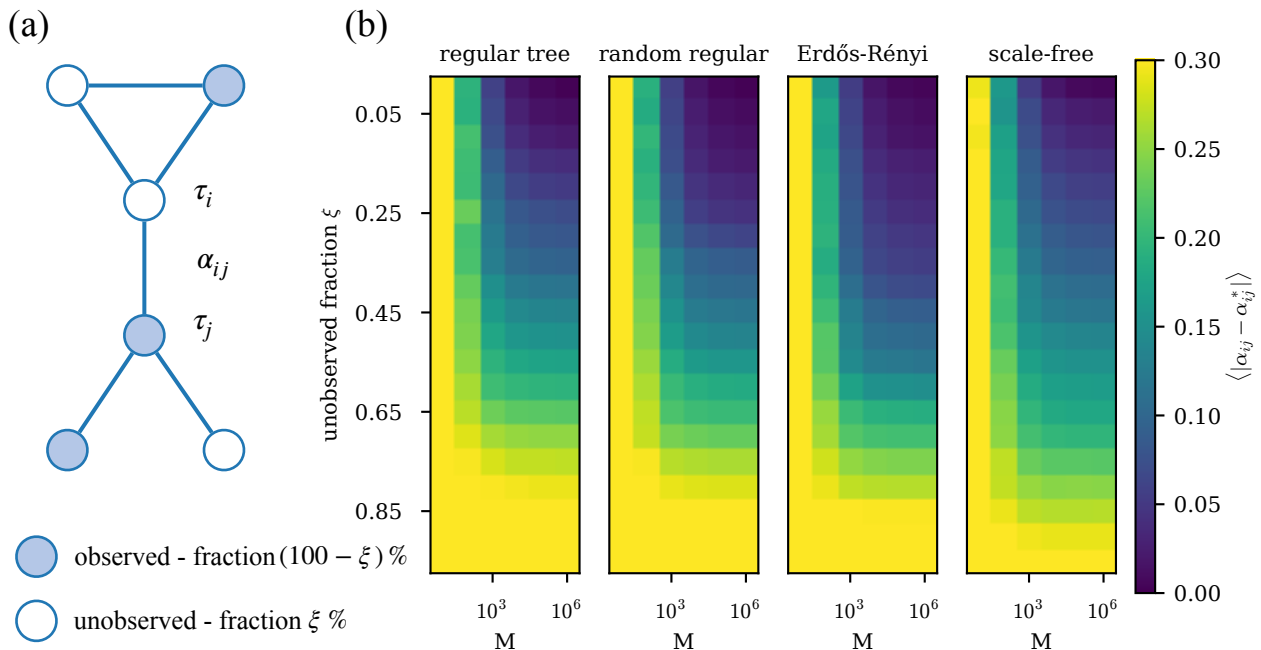
(a)

(b)



FIG. 2: Heat maps of average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades and the fraction of unobserved nodes. Each heat map represents a different network type, but all of them have the same color scale. Each point is averaged over 5 different networks and 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in the range $[0, 1]$). All networks contain $N = 100$ nodes and all but the tree have average degree $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length $T = 5$. Note that $\xi = 33\%$ corresponds to a random guess.

data presented in Figure 1. All subsections devoted to particular challenges are supported with simulations on synthetic networks. For this purpose we choose four different network models: 3-regular tree (RT), 3-regular random graph (RR), Erdös-Rényi graph with average degree equal to 3 (ER) [52] and a scale-free graph generated with Barabási-Albert model with average degree equal to 3 (BA) [53]. Except for the tree, all models have the same number of edges, so that they are comparable. The tree network case is studied for the baseline benchmarking purposes given that DMP is exact on trees. In the section III F, we apply the method to real-world networks, while assuming a simultaneous co-occurence of multiple types of incomplete data.

### A. SLICER's learning limits

SLICER was already shown in [8] to perform well in a regime where the unobserved part of the system is significant. It was not shown, however, what are the limits of the algorithm in terms of the size of the unobserved part and how it depends on the sample size. Here we answer this question and additionally we show how these limits depend on the network structure.

We take all four synthetic network models, generate up until $M = 10^6$ cascades and vary the percentage of the unobserved part from 0% to 95% (with a 5% step). The results of applying SLICER to all these cases are presented in the form of heat maps in Fig. 2. Interestingly, the quality of reconstructed parameters decrease faster for tree graphs than for loopy networks, because tree graphs have significantly lower edge density than other networks, making it not directly comparable. At the same time, single unobserved node affects a tree structure in a more significant way – single node separates a tree into two disconnected sub-graphs. We observe an interesting behavior from a scale-free graph. Despite a steeper decline in the quality of reconstruction for smaller fractions of the unobserved part $\xi$, for high values of $\xi$ the reconstruction still correlate with the true solution. If we look at $80 - 90\%$ of unobserved nodes for both random regular and Erdös-Rényi graph, the reconstructed parameters are basically random. Dynamics on scale free graph on the other hand, still does contain information that one is able to exploit, even for unobserved fraction equal to $\xi = 90\%$. For lower values of the unobserved fraction $\xi$, in the interval between $30 - 70\%$, the best parameters' reconstruction results are obtained for the Erdös-Rényi graph. Finally, when $\xi$ is below 30%, the best parameters' reconstruction is achieved for the regular random structure. These results illustrate that the uncertainty of the obtained reconstruction depends on the particular network structure.
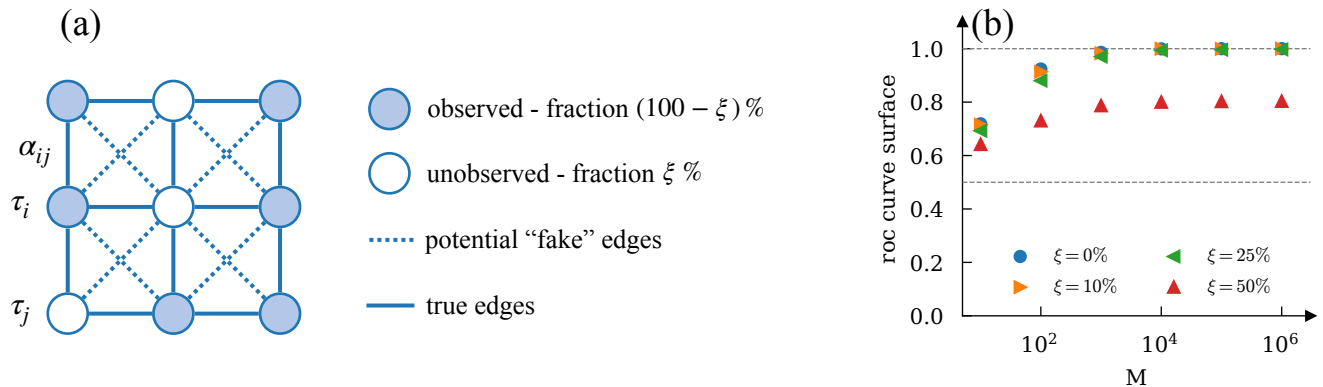
FIG. 3: (a) The structure of the lattice expanded by additional fake (inactive) edges. The structure learning task is to find the true edges, which were used to produce the observed cascades. (b) Average ROC curve surface, as a function of the number of available cascades for square lattice with additional fake (inactive) edges in the case where a fraction $\xi$ of nodes is unobserved. Each point is averaged over five different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in the range $[0, 1]$). Network contains $N = 100$ nodes. Unobserved nodes were picked at random. All cascades had length $T = 5$.
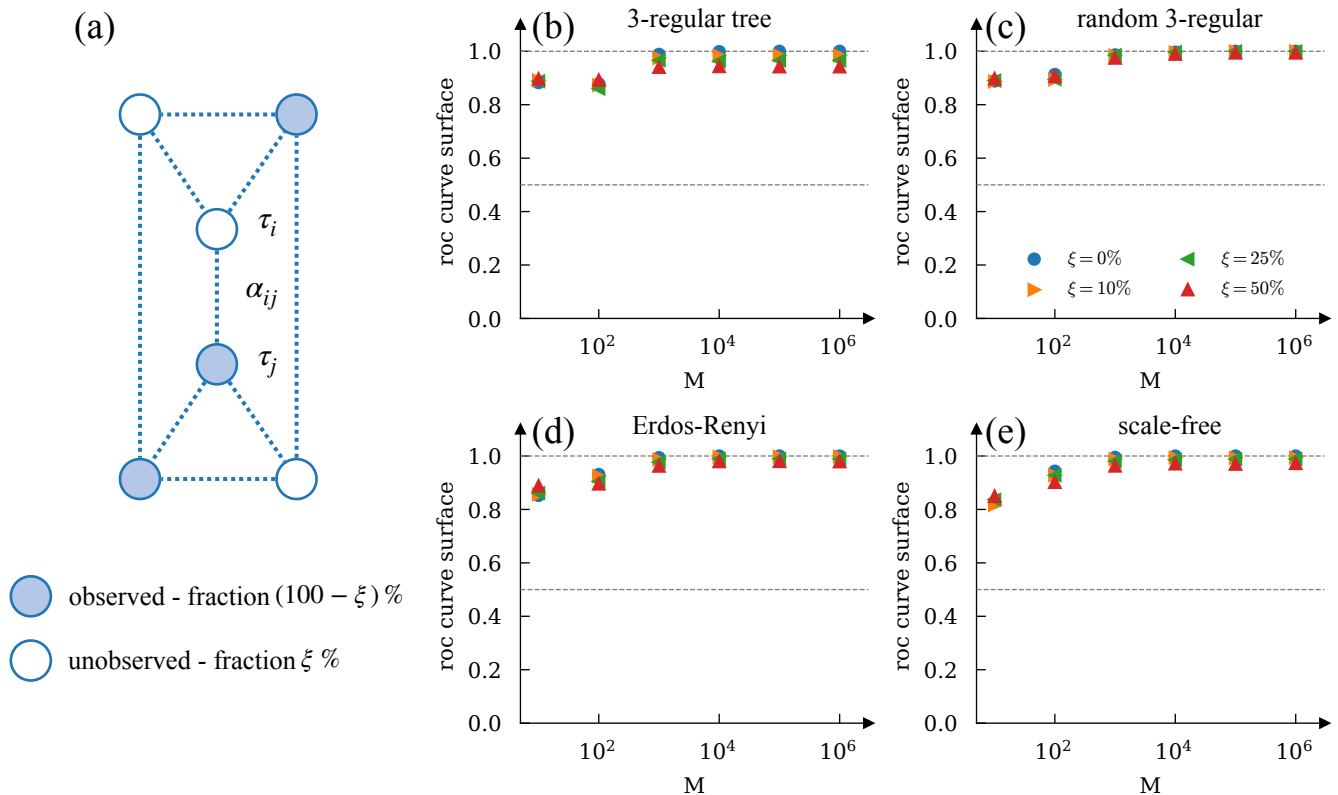


FIG. 4: Structure learning starting with a super-set of edges. Average ROC curve surface, as a function of the number of available cascades for different network types in the case where a fraction $\xi$ of nodes is unobserved and there is a known set of potential edges, where the number of *fake* edges is equal to the *true* ones. Each point is averaged over five different networks and five different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in the range $[0, 1]$). All networks contain $N = 100$ nodes and all but the tree have average degree $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length $T = 5$.

## B. Structure learning

In many real-world applications, like epidemic spreading, data on contact network is extremely difficult to gather. Sometimes it is possible to assess a limited number of potential connections, but detailed knowledge about the graph is rarely available. In practice, structure

learning is equivalent to learning spreading parameters, where near-zero couplings signal an absence of an edge. Under the structure recovery task, parameter learning needs to be run on larger graph reflecting the prior knowledge on the network structure. The less knowledge one have about the network, the larger the set of edges one needs to consider, up until $N(N-1)/2$ potential edges in a fully-connected graph when no knowledge about model structure is provided. When the dynamics of all nodes is available, edges can be recovered using maximum likelihood approach and thresholding. We focus on a more challenging scenario where the network structure is unknown or partially known under the presence of unobserved nodes. One of the main difficulty of this scenario consists in an increased potential for the solution degeneracy. Indeed, already in the case of full structural knowledge, degeneracy may appear for specific arrangements of unobserved nodes. The simplest example is when an unobserved node is a leaf – a node with a single connection. In this case it is impossible to recover the outgoing spreading parameter corresponding to the leaf edge. There are also more complex situations involving solution degeneracy such as interconnected clusters of hidden nodes, and increased number of unknown connections makes these situations more likely. However, information contained in the cascades makes it possible to discover the structure of the diffusion network with high accuracy even on very loopy graphs, as we show below.

Consider a problem of selecting the spreading graph from a set of known super-set of edges. Such a super-set of edges becomes a fully-connected graph in the worst case of no prior information on the diffusion network. SLICER can then be used as the structure discovery algorithm as follows. A gradient descent on the parameters $\{\alpha_{ij}\}$ as a part of SLICER is run on a super-set of edges. If a parameter goes beyond a certain threshold value (which we take as $10^{-8}$ in all experiments in this section), it is declared as zero and the respective edge is removed from the set of candidate edges, thus reducing the network at the next steps of the gradient descent procedure. This online pruning procedure results in a reduced overall computational complexity of the algorithm, compared to an alternative where SLICER is run on the full super-set of edges and the parameters are thresholded only in the end of the procedure. In order to evaluate the effectiveness of our method we use a receiver operating characteristic (ROC) curve [54]. After assessing all potential edges, we build the ROC curve of true and false positives among them and then we compute the surface under such curve. ROC curve surface equal to 1 represents a perfect reconstruction, while values oscillating around 0.5 suggest that the edge set selection is not better than a random guess.

Surprisingly, structure discovery with SLICER using the procedure described above is successful even when the underlying network has many short loops, which strongly affects the accuracy of the DMP approximation, as explained in section II B. To illustrate this point, consider an adversarial scenario of a regular two-dimensional square lattice containing a large number of short loops, which represents the ground-truth propagation network. Further, consider a super-set of edges by adding *diagonal* connections to the lattice, as shown in Fig. 3(a). Thus expanded network containing 50% fake edges (and even more short loops compared to the ground-truth square lattice) is then used as a starting network in the structure learning task. Ground-truth parameters $\alpha_{ij}$ are drawn from a uniform distribution over $[0, 1]$ interval. Despite this adversarial scenario, we are able to perfectly reconstruct the correct network, see Fig. 3, even with 25% of nodes being unobserved.

In Fig. 4, we test structure recovery procedure on four different types of synthetic networks, where the size of the super-set of edges is two times larger than the number of ground-truth edges in the network on which the data has been generated. These results show that an accurate recovery of the network structure is possible even in the presence of a large fraction of hidden nodes. In Section III F, we further explore structure learning for real-world network instances. In Appendix D, we discuss results for network structure recovery in the case of no available prior information on the topology, where the super-set of edges corresponds to a fully-connected graph.

### C. Simple graphs

In many applications, some knowledge about model parameters is assumed, which makes it easier to estimate them. Most often it is simply assumed that all the spreading parameters are the same for all edges $\alpha_{ij} = \alpha \, \forall_{(i,j) \in E}$. Although such setting significantly simplifies the problem, it is a good case study for understanding how the algorithm can take advantage of additional knowledge.

The constraint of equal parameters change the DMP equations, which leads to a different form of the SLICER+ instantiation for this case. In this new setting the constraints become a single parameter function:

$$\mathcal{L} = \underbrace{\mathcal{O}(\{\tau_i^c\})}_{objective} + \underbrace{\mathcal{C}(\alpha)}_{constrains} \,, \tag{9}$$

which requires recomputing the Lagrangian derivatives. Following the steps described in the Appendix B 2 we arrive at:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = -\frac{1}{\alpha} \sum_{s \in S} \sum_{t=0}^{T-1} \sum_{(i,j) \in E'} \lambda_{i \to j}^s(t) \cdot p_{i \to j}^s(t). \tag{10}$$

In the end, not only do we reduce the memory usage of the algorithm, but most importantly, we reduce the amount of data needed to obtain desirable level of error on parameters.

We apply the modified procedure to synthetic data generated with different network types. As shown in
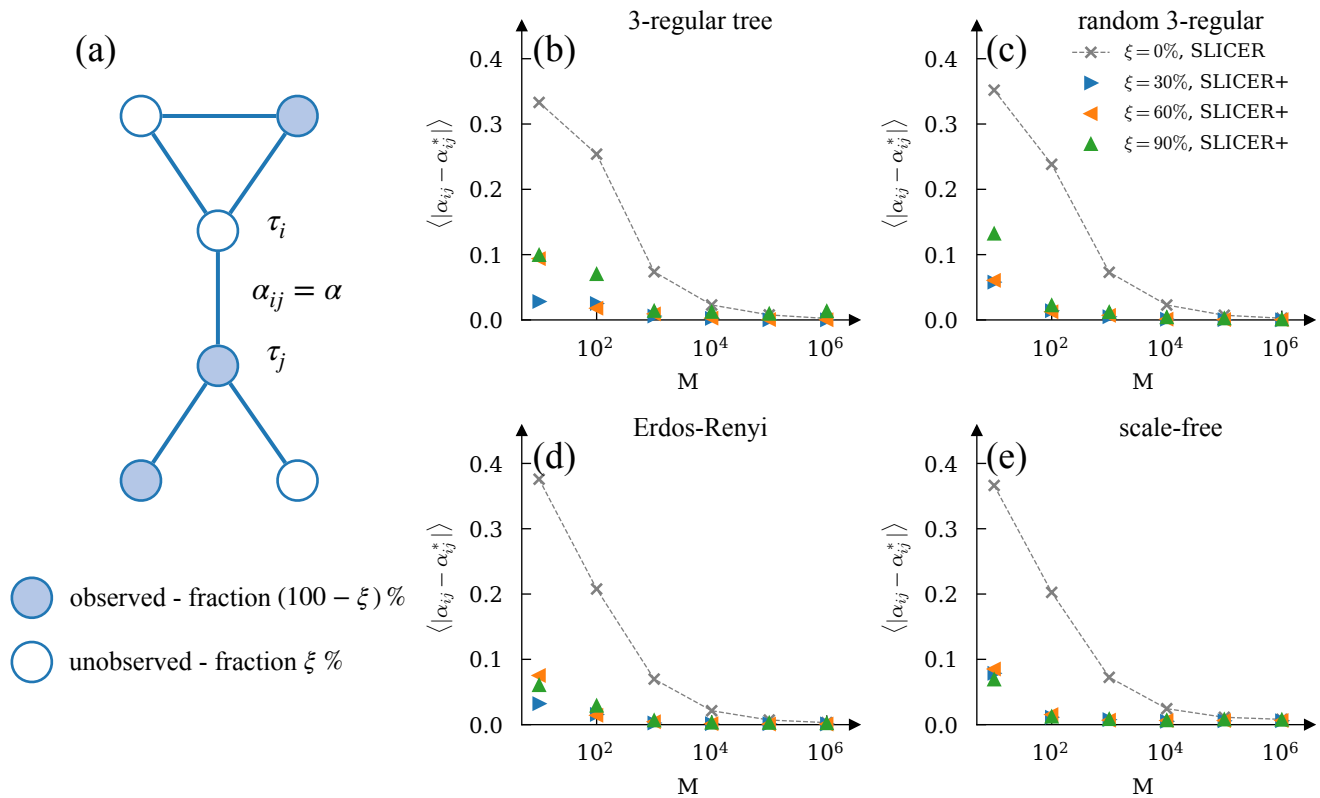
FIG. 5: Learning of simple graphs with equal transmission probabilities on all edges. Average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades for different network types in the case where a fraction $\xi$ of nodes is unobserved. The dashed gray line denotes the benchmark case where one is using standard SLICER without imposing the constraint of equal parameters on all edges. The triangles corresponds to the case where SLICER includes the knowledge about equal parameters. Each point is averaged over 5 different networks with parameters $\forall_{(i,j)\in E}\, \alpha_{ij}^* = 0.5$. All networks contain $N = 100$ nodes and all but the tree have average degree equal to $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length equal to $T = 5$.

Fig. 5 even with 90% unobserved nodes and regardless of the network structure, a small number of cascades is needed to obtain remarkable accuracy. Ten cascades are on average enough to get error below 0.1 and it does not depend strongly on the size of the unobserved part. Assuming different transmission probabilities for each edge in the reconstruction, shown with a gray dashed line as a benchmark, yields results which are orders of magnitude worse. Moreover, as shown in Fig. 2, different transmission probabilities combined with unobserved fraction above 60% yields almost random outcome, even with large number of cascades.

### D. Missing temporal observations

Not being able to observe the whole system, as a result of a limited budget or other challenges related to gathering data, is a fair limitation for many practical applications. Another realistic situation arises when, for similar reasons, one is not able to record the data for ex-

tended periods of time. Partial observability of the system, should be considered in both spatial and temporal spaces. Here we show how to modify the objective so that it correctly accounts for the unobserved time periods, resulting in a SLICER+ instantiation for the scenario of missing observations in times.

Under this scenario, the state of the network is not observed at each time step, but only at a subset of observation times. Therefore, although we do not generally observe activation time of a node, we still know that activation happened somewhere in the last unobserved period preceding the first time we noticed a given node to be active. As a result there are two cases: (i) node activation is observed and we denote the activation time as $\tau_i$; (ii) node activation happens during an unobserved interval denoted as $[\hat{\tau}_i - \Delta\hat{\tau}_i, \hat{\tau}_i]$. Observed snapshots of the initial and of the final states of the spreading dynamics only represents a particular example of such a missing information in time. New objective will depend on known
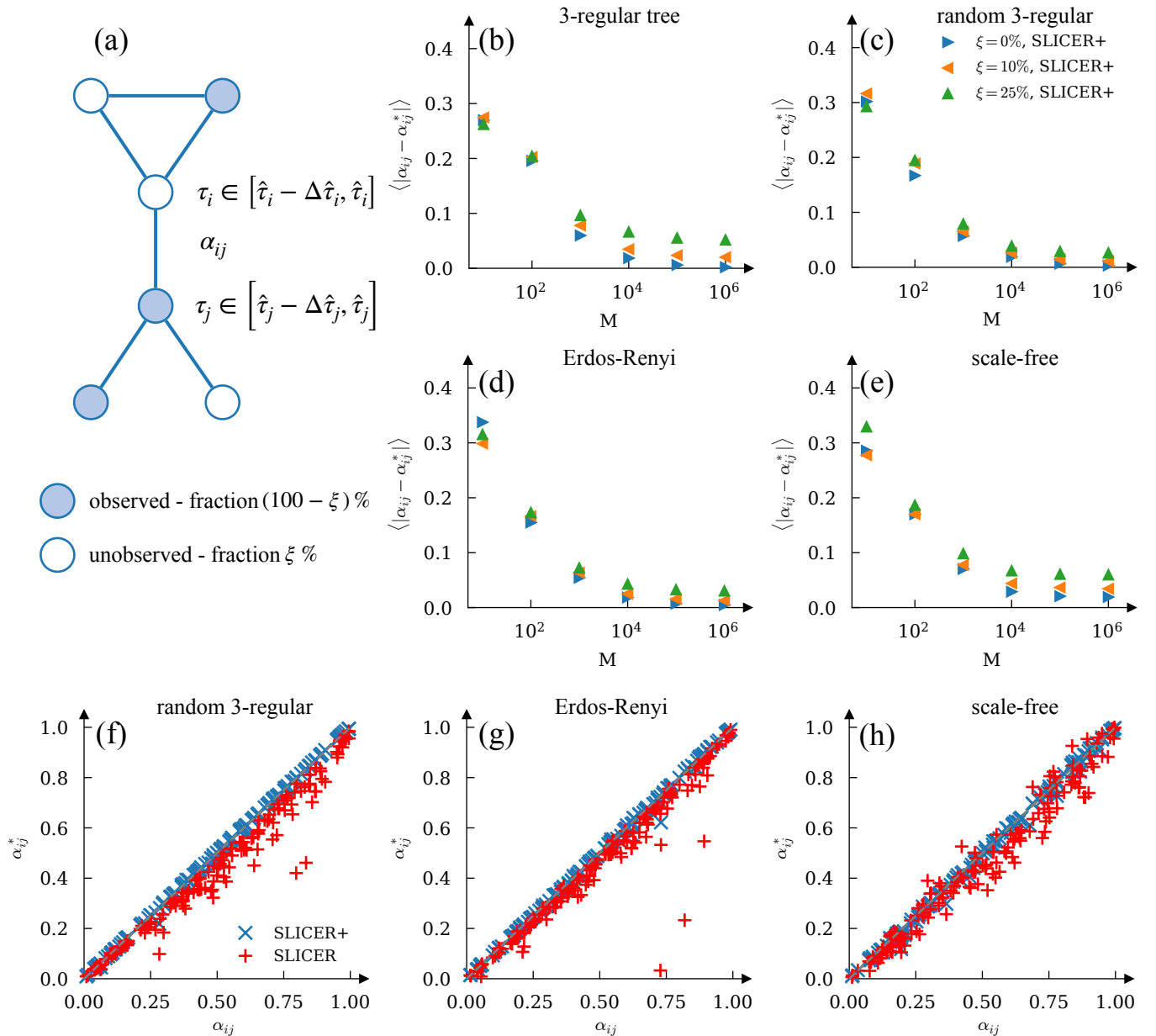
FIG. 6: **(a)** Reconstruction of parameters under missing temporal observations. **(b-e)** Average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades for different network types. Only initial and final spreading states are observed and additionally, some nodes are not observed at all. The triangles corresponds to the case where SLICER includes the knowledge about unobserved periods of time. Each point is averaged over 5 different networks and 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in $[0, 1]$). All networks contain $N = 100$ nodes and all but the tree have average degree equal to $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length equal to $T = 6$. **(f-h)** Scatter plots of true parameters $\alpha_{ij}^*$ versus the estimated ones $\alpha_{ij}$. The results were obtained for different network structures with $N = 100$ nodes, $M = 10^6$ cascades of length $T = 6$ and $\xi^{time} = 33.(3)\%$ of unobserved times picked at random. Except for the tree, all networks have average degree equal to $\langle k \rangle = 3$.

activation times and activation intervals:

$$\mathcal{L} = \underbrace{\mathcal{O}(\{\tau_i^c\}, \{\hat{\tau}_i, \Delta\hat{\tau}_i\})}_{objective} + \underbrace{\mathcal{C}(\{\alpha_{ij}\})}_{constrains}. \quad (11)$$

This leads to a variant of the SLICER+ algorithm for

missing temporal observations, presented in Appendix B 3.

Results in Fig. 6(b-e) are presented in the case where only initial and final states of the dynamics are observed, similarly to the case considered in [41] (where the case of
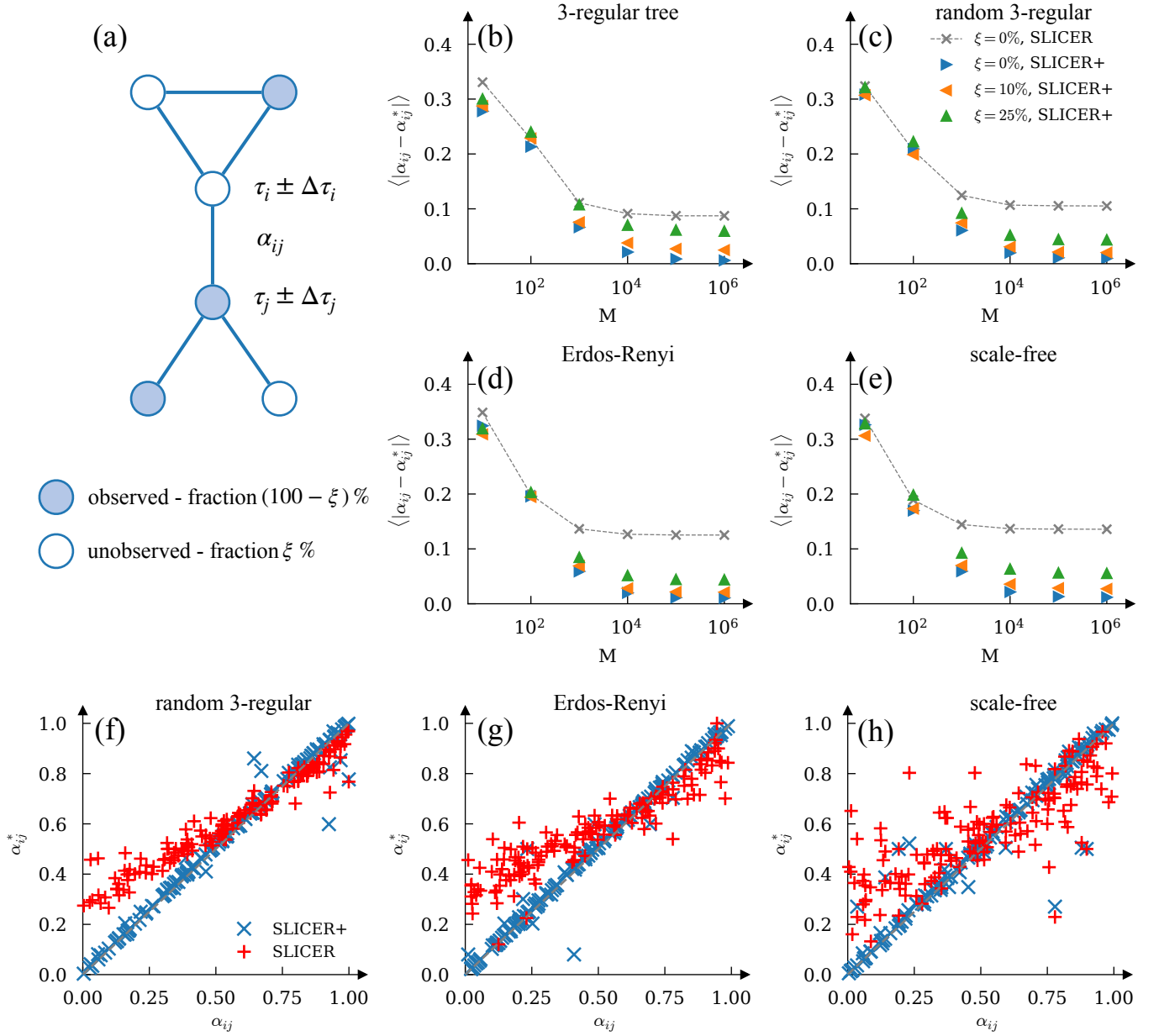
FIG. 7: **(a)** Reconstruction of parameters under noisy observations. **(b-e)** Average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades for different network types in the case where a fraction $\xi$ of nodes is unobserved. Additionally, observed activation times $\{\tau_i^c\}_{i \in O}$ are subject to noise described by Eq. 14. The dashed gray line denotes the benchmark case obtained with SLICER. The triangles corresponds to the case where SLICER includes the knowledge about the noise. Each point is averaged over 5 different networks and 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in $[0, 1]$). All networks contain $N = 100$ nodes and all but the tree have average degree equal to $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length equal to $T = 5$. **(f-h)** Scatter plots of true parameters $\alpha_{ij}^*$ versus the estimated ones $\alpha_{ij}$. The results were obtained for different network structures with $N = 100$ nodes, $M = 10^6$ cascades of length $T = 5$ and $\xi = 10\%$ of unobserved nodes. Additionally, observed activation times $\{\tau_i^c\}_{i \in O}$ are subject to noise described by Eq. 14. Except for the tree, all networks have average degree equal to $\langle k \rangle = 3$.

full spatial observations $\xi = 0$ has been treated). Given the effectiveness of the method in the setting of unobserved time-periods, we explore an even more complex case, where we combine both unobserved times and un-

observed nodes, a realistic setting, which, to our best knowledge, was not addressed in the literature before. Results presented in Fig. 6(b-e) are of similar quality as the ones obtained in [8] when there are no missing ob-

servations in time, despite lack of knowledge on the exact activation times, which illustrates the robustness of SLICER+ with respect to missing observations in time.

Unobserved time periods can be naively treated with SLICER, by simply assuming that if a given node was activated in such a period, it is regarded as an unobserved node. This, however, does not take advantage of all the available information. As it is shown in Fig. 6(f-h), this leads to a one-sided bias, which quantitatively depends on the network structure and the fraction of unobserved timestamps. When the unobserved periods are treated correctly, as done in SLICER+, the bias disappears.

### E. Noisy timestamps

Another realistic scenario in data collection relates to the fact that apart from not being able to monitor everything, the collected data is not necessary perfect. This may be a result of omission, deliberate actions or simply imperfect gathering procedure. If, for example, a population is tested for a certain virus, the testing dates are most likely not the dates when the infection occurred. On top of that, the lab results may be delayed, dates may be mistaken and tests can produce both false positives and false negatives. All of this builds up into an unpredictable noise, which cannot be removed in a deterministic matter. It has to be accounted for in the modelling process.

In our case of learning spreading parameters, the data is a set of activation times. In this context, we model noisy observations as deviations around the correct activation times. We assume that what is observed is a sum of the original activation time and a noise coming from a certain discrete distribution. In general it is described as follows:

$$P(\tau = \tau^* + k) = \pi_k \quad \forall |k| \leq K, \qquad (12)$$

where $\tau$ is the observed activation time, $\tau^*$ is the real activation time, and $K$ is a constant related to the support of the noise distribution. The presence of noise in timestamps has no effect on the DMP equations and the constraints, but will affect the objective function, resulting in incorporation of the prior on the noise in the SLICER+ formulation. The objective function will now depend on both activation times and the noise distribution:

$$\mathcal{L} = \underbrace{\mathcal{O}(\{\tau_i^c\}, \{\pi_k\})}_{objective} + \underbrace{\mathcal{C}(\{\alpha_{ij}\})}_{constrains}. \qquad (13)$$

Exact form of the objective for noisy data and detailed derivation of SLICER+ for this case is given in the Appendix B 4.

As in the previous sections, we test the instantiation of SLICER+ for the case of noisy observations with different network structures. In the numerical experiments we use noise described by the following distribution for

$K = 1$:

$$\pi_k = \begin{cases} 0.6, & \text{if } k = 0, \\ 0.2, & \text{if } |k| = 1, \\ 0.0, & \text{otherwise.} \end{cases} \qquad (14)$$

For simplicity, in the numerical experiments in this section, we assume that the noise distribution is known, but the probabilities $\pi_k$ characterizing the noise distribution can be treated as parameters and learned from the data.

Notice that a naive application of SLICER to the case of noisy observations could lead to infinite objective function, since the realization of noisy timestamps could be *impossible*, showing inconsistency with the dynamic rules of the IC model. This could happen, for instance, if a node is observed to be activated before any of its neighbor is activated. Application of SLICER is still possible by neglecting these conflicting cases leading to an infinite value of the objective function. We use such an approach as a benchmark to SLICER+ that accounts for the possibility of noisy observations. Fig. 7(b-e) shows that the use of SLICER quickly leads to a significant error gap, while consideration of noise implemented in SLICER+ brings the algorithm's performance to a level similar to a noiseless situation. Looking more closely at the particular results at Fig. 7(f-h), we see that disregarding the presence of noise gives results that are correlated with the true solution, but there is a systematic bias, which depends on the network structure. The solution given by the noise-corrected algorithm, on the other hand, aligns perfectly with the ground-truth.

### F. Study on real-world networks

Previously we used different network structures to show how they affect the results of the learning procedure for different scenarios. Real-world networks typically combine several characteristics of different random models [55]. In order to truly test the effectiveness of our learning framework we use two real-world social networks together with multiple types of noisy and incomplete spreading data. The first network is the Zachary Karate Club [56] with $N = 34$ nodes, $|E| = 78$ edges and multiple short loops, small enough to serve as an explicit illustration of the quaility of structure learning under partial and noisy observations. The second network is a Facebook snapshot from [57] and previously analysed in [58], with $N = 2888$ nodes and $|E| = 2981$ edges.

In Figure 8, we present results on learning of the structure of the network from a combination of missing information: partially reporting nodes and noisy activation times. We assume that no prior information on the topology is known to the algorithm. Due to many short loops, learning spreading parameters in the Zachary Club network using SLICER+ is not an easy task. As shown in Fig. 8(c), 10% of unobserved nodes leads to a low error in parameter recovery of 5% error on average. With 25% of
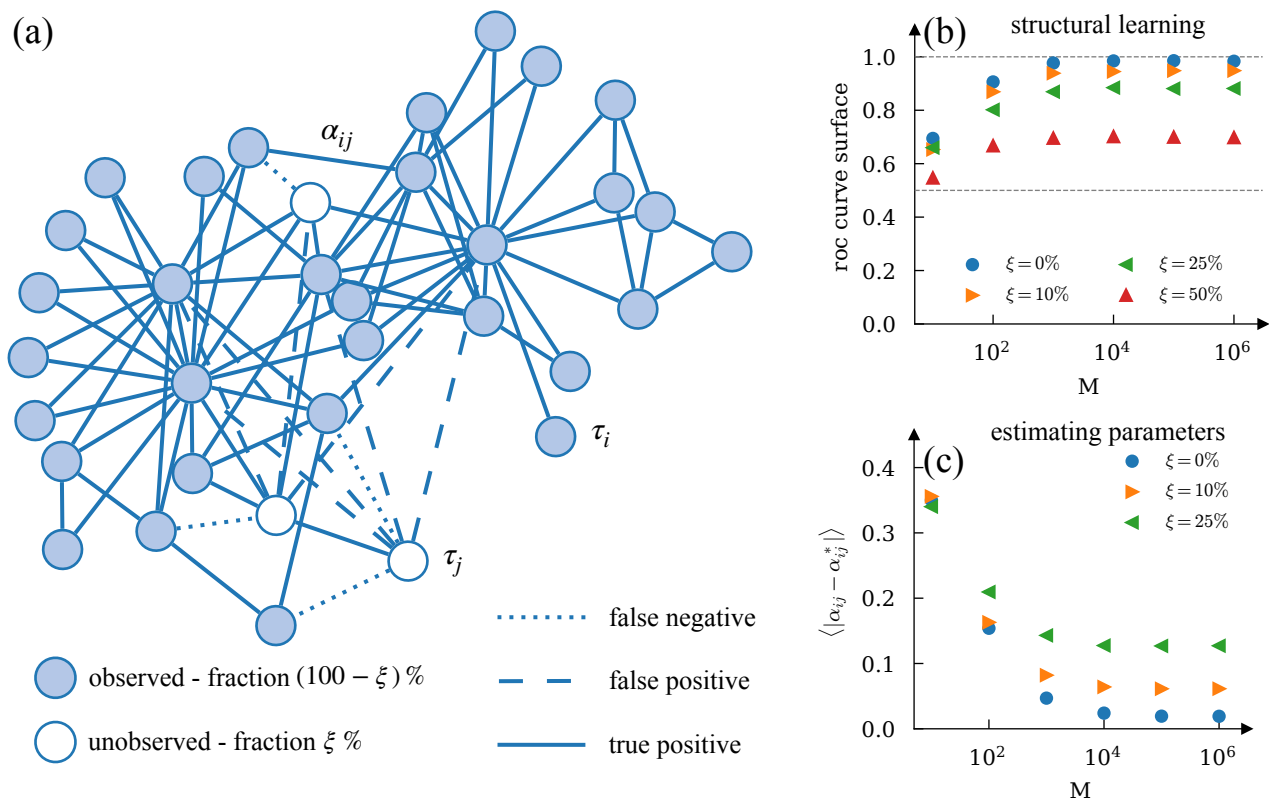
FIG. 8: (a) Structure reconstruction using the Zachary Karate Club network. Solid black connections represent the correctly reconstructed edges. Dotted black connections are the ground-truth edges that were not reconstructed. Dashed gray lines are the ones reconstructed by the algorithm, but absent in the ground-truth edge set. We fix all of spreading parameters to $\alpha_{ij} = 0.5$. Network was learned based on $M = 1000$ cascades of length $T = 5$ with $\xi = 10\%$ of unobserved nodes. (b) Average ROC curve surface, as a function of the number of available cascades for the Zachary Karate Club network in the case where a fraction $\xi$ of nodes is unobserved. (c) Average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades for the Zachary Karate Club network in the case where a fraction $\xi$ of nodes is unobserved. In both cases (b) and (c), each point is averaged over 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in $[0, 1]$). Network contains $N = 34$ nodes and $|E| = 78$ edges. Unobserved nodes were picked at random. All cascades had length $T = 5$. We assume no knowledge on the topology, hence the initial set of edges is the complete graph. Additionally, observed activation times are subject to known noise described in Eq. 14.

unobserved nodes, this error grows above 10%. However, the structure recovery is performed with a greater accuracy, despite no assumed knowledge on the edge set. Obviously, no algorithm can produce perfect reconstruction for high fractions of non-reporting nodes and unobserved edges due to a degeneracy in the space of solutions.

An explicit structure learning example produced for $\xi = 10\%$ unobserved nodes and $M = 1000$ cascades of length $T = 5$ is shown in Fig. 8(a). To eliminate possibility of bias due to a choice of ground-truth $\alpha_{ij}$, we use the test case with all $\alpha_{ij} = 0.5$. The algorithm however does not have access to this information during the learning procedure. Although the parameter reconstruction quality could be slightly improved by increasing the number of cascades, one can already see that the reconstruction on the observed part is perfect. Incorrect classification appears only in the vicinity of the unobserved

nodes. Specifically, the problem is more pronounced for the two connected unobserved nodes, whereas the result is much better for an isolated unobserved node.

We now test structure recovery with SLICER+ on the larger Facebook network with the following setting: the observed timestamps are noisy, and we start with a superset of 5962 edges, twice as large as the ground-truth set of 2981 edges. The Facebook snapshot network is a sparse graph with multiple hubs. Hubs tend to increase the number of triangles, but the clustering of this network is on average lower than for the Zachary Club. This is consistent with our observation that the reconstruction error in spreading parameters is significantly smaller, even when 25% of nodes are unobserved, as shown in Fig. 9. Furthermore, inclusion of the knowledge on the super-set of 5962 potential edges allows to nearly perfectly reconstruct the network structure in the partial observation

regime, assuming enough data is available.

## IV. CONCLUSIONS

Literature on learning of spreading models from data is rich, but not much attention so far was given to the problem of learning from incomplete and uncertain information. Our work proposes a flexible method that addresses this problem. SLICER+ deals with multiple settings related to incomplete and noisy data, that include partial observations in both spatial and temporal dimensions, noisy timestamps, unknown network structure, and combinations of these settings.

For unobserved times, we show that a proper formulation of the algorithm is able to learn model parameters from final and initial states only, assuming that enough data available. This holds even in the presence of unobserved nodes. We further show that uncertainty in provided activation times leads to significant bias, when not properly accounted for.

Although the structure learning task under unobserved nodes may generate degenerate solutions, we empirically showed that additional knowledge can significantly increase the quality of reconstruction using several representative synthetic networks, as well as two popular real-world networks. For real-world instances specifically, we used a setting where different types of incomplete and uncertain data settings are combined.

We found that our results are quite robust regardless of the synthetic network type. As an illustration, we considered quality of reconstruction as a function of the size of the set of unobserved nodes on a variety of random graph families.

Our framework can be further developed in several future directions. In the setting of noisy observations, we assumed for simplicity that the noise distribution is known. However, the noise distribution probabilities $\pi_k$ can be treated as additional parameters and learned from the data. For the structure learning task, we considered reconstruction from a super-set of edges. It would be interesting to test the impact of other graph topological prior information on learning, such as details on density or degree distribution. Such a setting would be relevant in the setting where surveillance and data gathering procedures can be controlled. Under the scenario where the observations are too sparse and the degeneracy in the reconstruction is unavoidable and in applications focused on downstream prediction tasks, it could be useful to study the properties of the learned models as *effective* representations [8]. Finally, in the future, it would be useful to extend our framework to other dynamic models, including models with reversible dynamics.

### CODE AND SUPPLEMENTARY MATERIAL

Full implementation of all algorithms studied in this work is available at [59].

### ACKNOWLEDGMENTS

### Appendix A: Derivation of the DMP equations for the IC model

The DMP equations presented in the paper are quite intuitive. In case of simple models, such as IC model, they could be derived by identifying the correct dynamical variables to use in the equations, as it was done here. They are also intrinsically connected to a more general framework of Belief Propagation (BP) equations, which are well described in detail in the literature [29, 30]. In the dynamic setting, the BP equations on time trajectories have been studied by Kanoria and Montanari in [11]. Here, for completeness, we show how DMP equations for IC model can be obtained directly from BP formulation.

Denote $\vec{\sigma}_i = (\sigma_i^0, \ldots, \sigma_i^T)$ the time trajectory of node $i$ variable, and $w_i\left(\sigma_i^{t+1} | \sigma_i^t, \{\sigma_j^t\}_{j \in \partial i}\right)$ the local probability of node $i$ transitioning to state $\sigma_i^{t+1}$ at time $t+1$, given its previous state and the states of its neighborhood at time $t$, and $P(\{\sigma_i^0\}_{i \in V})$ is the probability of the initial state. Additionally, let us introduce new type of message $\mu^{i \to j}(\vec{\sigma}_i || \vec{\sigma}_j)$, which represents the probability of node $i$ variable having the trajectory $\vec{\sigma}_i$ on an auxiliary graph, where the trajectory of node $j$ is fixed to $\vec{\sigma}_j$. Our starting point is the dynamic belief propagation (DBP) equation on time trajectories [11]:

$$\mu_t^{i \to j}(\vec{\sigma}_i || \vec{\sigma}_j) = \sum_{\{\sigma_k^0, \ldots \sigma_k^{t-1}\}_{k \in \partial i \setminus j}} P(\{\sigma_i^0\}_{i \in V}) \prod_{\tau=0}^{t-1} w_i\left(\sigma_i^{\tau+1} | \sigma_i^\tau, \{\sigma_k^\tau\}_{k \in \partial i}\right) \prod_{k \in \partial i \setminus j} \mu_{t-1}^{k \to i}(\vec{\sigma}_k || \vec{\sigma}_i), \tag{A1}$$

where $\mu_t^{i \to j}(\vec{\sigma}_i || \vec{\sigma}_j)$ is a message probability for trajectories up until time $t$.

A direct use of the above equations is still impractical due to their exponential complexity over cascade length $T$. This is where the unidirectional character of the dynamics come in handy and leads to simplifications. In the case
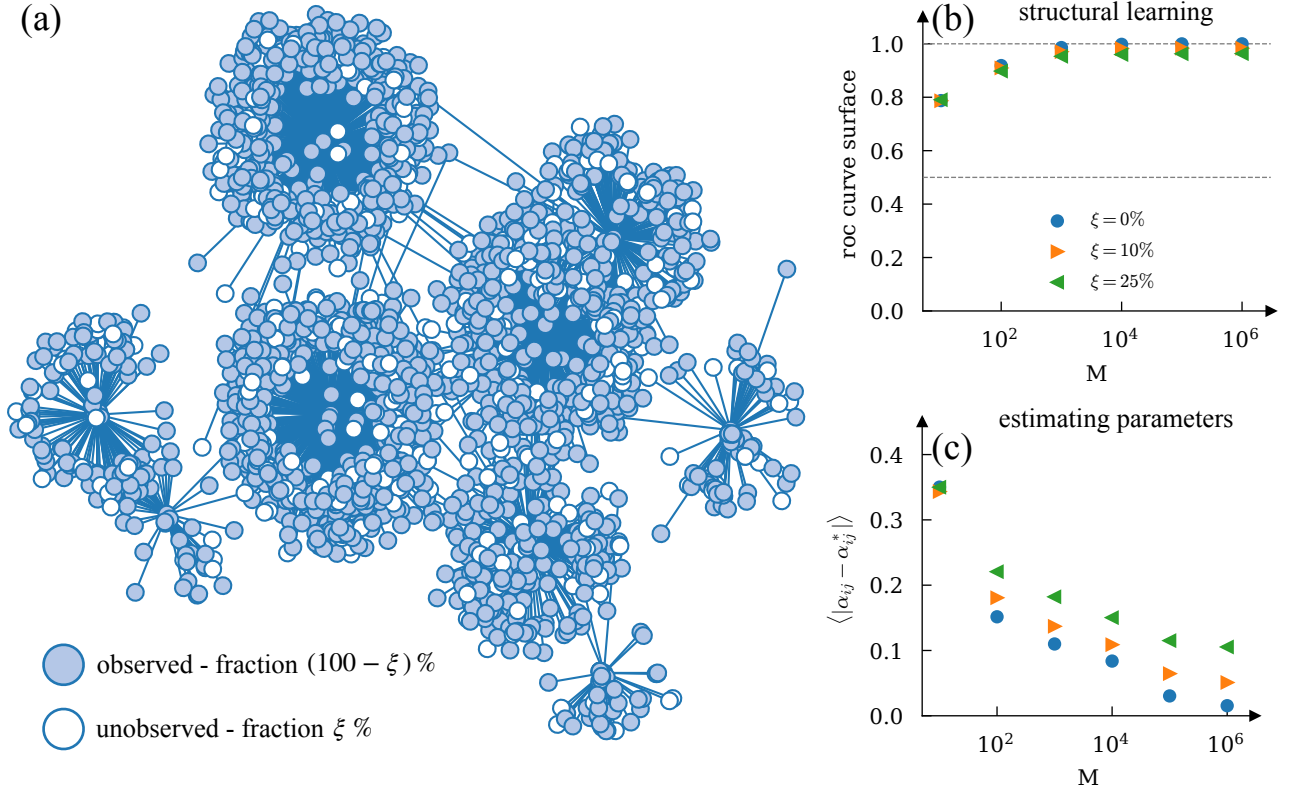
FIG. 9: (a) Structure of the FB network. (b) Average ROC curve surface, as a function of the number of available cascades for the Facebook network in the case where a fraction $\xi$ of nodes is unobserved. (c) Average difference between inferred $\alpha_{ij}$ and real $\alpha_{ij}^*$ parameters in $\ell_1$ norm, as a function of the number of available cascades for the Facebook network in the case where a fraction $\xi$ of nodes is unobserved. In both cases, each point is averaged over 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in $[0, 1]$). Network contains $N = 2888$ nodes and $|E| = 2981$ edges. Unobserved nodes were picked at random. All cascades had length $T = 5$. We assume that a super-set of 5962 edges is known, out of which we would like to infer the ground-truth set of 2981 edges. Additionally, observed activation times are subject to known noise described in Eq. 14.

of the IC model, vector $\vec{\sigma}_i$ can be fully described by a single number $\tau_i$, which represents the activation time. This allows us to rewrite the DBP equations as follows:

$$\mu_{T+1}^{i \to j}(\tau_i \| \tau_j) = \sum_{\{\tau_k\}_{k \in \partial i \setminus j}} \bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) \prod_{k \in \partial i \setminus j} \mu_T^{k \to i}(\tau_k \| \tau_i), \tag{A2}$$

$$\mu_{T+1}^i(\tau_i) = \sum_{\{\tau_k\}_{k \in \partial i}} \bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) \prod_{k \in \partial i} \mu_T^{k \to i}(\tau_k \| \tau_i), \tag{A3}$$

where $\mu_T^{i \to j}(\tau_i \| \tau_j)$ and $\mu_T^i(\tau_i)$ correspond the same messages and marginals as before, but parameterized with activation times $\tau_i$. Additionally, we denote $\bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) = P(\{\sigma_i^0\}_{i \in V}) \cdot W_i(\tau_i; \{\tau_k\}_{k \in \partial i})$ for simplicity. Before we move further we want to highlight some observations, which will be used later. First, all the information required to compute the probability of activation at certain time $\tau_i$ is available at that time.

$$\mu_T^i(\tau_i) = \mu_{\tau_i}^i(\tau_i) \quad \forall_{T > \tau_i}, \tag{A4}$$

$$\mu_T^{i \to j}(\tau_i \| \tau_j) = \mu_{\tau_i}^{i \to j}(\tau_i \| \tau_j) \quad \forall_{T > \tau_i}. \tag{A5}$$

Second, if the activation time $\tau_i$ of the node $i$ is earlier then the activation time $\tau_k$ of its neighbor $k$, than the message from $i$ to $k$ does not depend on $\tau_k$. If we combine this we the previous equation, we can write:

$$\mu_T^{i \to k}(\tau_i \| \tau_k) = \mu_T^i(\tau_i \| \infty) \quad \forall_{\tau_k > T}, \tag{A6}$$

where $\infty$ describes the state where a given node was activated at any time later than the cutoff time $T$ (the time subscript). Finally, the incoming messages in equations above are independent of node $i$ activation time $\tau_i$ and always behave as $\tau_i$ is later than their cutoff time:

$$\mu_T^{i \to j}(\tau_i || \tau_j) = \sum_{\{\tau_k\}_{k \in \partial i \setminus j}} \bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) \prod_{k \in \partial i \setminus j} \mu_{T-1}^{k \to i}(\tau_k || \infty), \tag{A7}$$

$$\mu_T^i(\tau_i) = \sum_{\{\tau_k\}_{k \in \partial i}} \bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) \prod_{k \in \partial i} \mu_{T-1}^{k \to i}(\tau_k || \infty). \tag{A8}$$

This is a direct consequence of the previous observations. Let us now define new type of marginals and messages:

$$P_S^i(t) = \sum_{\tau_i > t} \mu_T^i(\tau_i) = \mu_t^i(*), \tag{A9}$$

$$P_S^{i \to j}(t) = \sum_{\tau_i > t} \mu_T^{i \to j}(\tau_i || \infty) = \mu_t^{i \to j}(* || \infty), \tag{A10}$$

where $P_S^i(t)$ is the probability that node $i$ was not activated until time $t$ and $P_S^{i \to j}(t)$ is the same probability, but on an auxiliary graph where node $j$ does not exist. Symbol $*$ stands for all the trajectories with the activation time $\tau_i$ being after the time horizon $T$. It means that when we sum over $\tau_i > b$, in practice we sum over $\tau_i \in \{b+1, b+2, \ldots, T-1, T, *\}$. It should be noted that in all equations below we use the convention that when $\tau_i = *$, $\tau_i - 1 = T$. Now, we use these new objects to compute the term $\bar{W}_i(\ldots)$, which describes the dynamics of the process. In the case of IC model it has the following form:

$$\begin{aligned}
\bar{W}_i(\tau_i; \{\tau_k\}_{k \in \partial i}) =\ & (1 - P_S^i(0)) \cdot \mathbb{1}(\tau_i = 0) \\
& + P_S^i(0) \cdot \mathbb{1}(\tau_i > 0) \left[ \prod_{k \in \partial i} 1 - \alpha_{ki} \cdot \mathbb{1}(\tau_k < \tau_i - 1) \right] \\
& \times \left[ 1 - \prod_{k \in \partial i} (1 - \alpha_{ki} \cdot \mathbb{1}(\tau_k = \tau_i - 1)) \cdot \mathbb{1}(\tau_i \neq *) \right]
\end{aligned} \tag{A11}$$

We can now plug it all together and get:

$$\begin{aligned}
1 - P_S^{i \to j}(t) &= 1 - \sum_{\{\tau_k\}_{k \in \partial i \setminus j}} \bar{W}_i(*; \{\tau_k\}_{k \in \partial i}) \prod_{k \in \partial i \setminus j} \mu_{t-1}^{k \to i}(\tau_k || \infty) \\
&= 1 - P_S^i(0) \sum_{\{\tau_k\}_{k \in \partial i \setminus j}} \prod_{k \in \partial i} (1 - \alpha_{ki} \cdot \mathbb{1}(\tau_k \neq *)) \prod_{k \in \partial i \setminus j} \mu_{t-1}^{k \to i}(\tau_k || \infty) \\
&= 1 - P_S^i(0) \sum_{\{\tau_k\}_{k \in \partial i \setminus j}} \prod_{k \in \partial i \setminus j} (1 - \alpha_{ki} \cdot \mathbb{1}(\tau_k \neq *)) \, \mu_{t-1}^{k \to i}(\tau_k || \infty) \\
&= 1 - P_S^i(0) \prod_{k \in \partial i \setminus j} \sum_{\tau_k} (1 - \alpha_{ki} \cdot \mathbb{1}(\tau_k \neq *)) \, \mu_{t-1}^{k \to i}(\tau_k || \infty) \\
&= 1 - P_S^i(0) \prod_{k \in \partial i \setminus j} \left( \sum_{\tau_k \leq t-1} (1 - \alpha_{ki}) \mu_{t-1}^{k \to i}(\tau_k || \infty) + \mu_{t-1}^{k \to i}(* || \infty) \right) \\
&= 1 - P_S^i(0) \prod_{k \in \partial i \setminus j} \left( 1 - \alpha_{ki}(1 - P_S^{k \to i}(t-1)) \right),
\end{aligned} \tag{A12}$$

which is equivalent to the DMP equations (3) from the main text, where $p_{i \to j}(t) = 1 - P_S^{i \to j}(t)$ and assuming some initial condition $\bar{p}_{i \to j} = 1 - P_S^{i \to j}(0)$. Readers interested in derivation of DMP from DBP for different dynamic models should refer to [11, 17, 23].

## Appendix B: Derivation of SLICER+ for different scenarios of incomplete data

Below we present the details of Lagrangian formulation used in the optimisation scheme of each of the different case of missing or incomplete data. We show step-by-step how to compute all the Lagrange multipliers and the gradient step used to update the parameters.

### 1. SLICER

The SLICER algorithm [8] was developed to deal with the scenario of missing information on nodes exclusively. In this setting the objective takes the following form:

$$
\mathcal{O} = \sum_{s \in S} \sum_{i \in O} \sum_{\tau_i^s} m^{\tau_i^s} \log \left( p_i^s(t) \cdot \mathbb{1}_{(t \leq T)} - p_i^s(t-1) \cdot \mathbb{1}_{(t>0)} + \mathbb{1}_{(t=*)} \right). \tag{B1}
$$

The constraints are the same as in the main text, but lets remind them for reader's convenience.

$$
\begin{aligned}
\mathcal{C} = &\sum_{s \in S} \sum_{t=0}^{T-1} \sum_{i \in V} \lambda_i^s(t+1) \left( p_i^s(t+1) - 1 + \left( 1 - \bar{p}_i^s \right) \prod_{k \in \partial i} \left( 1 - \alpha_{ki} \cdot p_{k \rightarrow i}^s(t) \right) \right) \\
&+ \sum_{s \in S} \sum_{t=0}^{T-1} \sum_{(i,j) \in E} \lambda_{i \rightarrow j}^s(t+1) \left( p_{i \rightarrow j}^s(t+1) - 1 + \left( 1 - \bar{p}_i^s \right) \prod_{k \in \partial i \backslash j} \left( 1 - \alpha_{ki} \cdot p_{k \rightarrow i}^s(t) \right) \right).
\end{aligned} \tag{B2}
$$

The sum of the objective and the constraints constitutes the Lagrangian. Lagrange multipliers can now be found by differentiating the Lagrangian:

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial p_i^s(t)} = &\lambda_i^s(t) + \sum_{\tau_i^s} \frac{m^{\tau_i^s} \cdot \mathbb{1}_{(t=\tau_i^s)} \cdot \mathbb{1}_{(\tau_i^s \leq T)}}{p_i^s(\tau_i^s) - p_i^s(\tau_i^s - 1) \cdot \mathbb{1}_{(\tau_i^s > 0)}} \\
&+ \sum_{\tau_i^s} \frac{m^{\tau_i^s} \cdot \mathbb{1}_{(t=\tau_i^s-1)} \cdot \mathbb{1}_{(\tau_i^s > 0)}}{p_i^s(\tau_i^s - 1) - p_i^s(\tau_i^s) \cdot \mathbb{1}_{(\tau_i^s \leq T)} - \mathbb{1}_{(\tau_i^s=*)}},
\end{aligned} \tag{B3}
$$

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial p_{i \rightarrow j}^s(t)} = &\lambda_{i \rightarrow j}^s(t) - \lambda_j^s(t+1)\, \alpha_{ij} \left( 1 - \bar{p}_j^s \right) \prod_{m \in \partial j \backslash i} \left( 1 - \alpha_{mj} \cdot p_{m \rightarrow j}^s(t) \right) \\
&- \sum_{k \in \partial j \backslash i} \lambda_{j \rightarrow k}^s(t+1)\, \alpha_{ij} \left( 1 - \bar{p}_j^s \right) \prod_{m \in \partial j \backslash \{i,k\}} \left( 1 - \alpha_{mj} \cdot p_{m \rightarrow j}^s(t) \right)
\end{aligned} \tag{B4}
$$

and equating the derivatives to zero. Finally, derivatives over parameters $\alpha_{ij}$ read

$$
\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \alpha_{ij}} = &-\sum_{s \in S} \sum_{t=0}^{T-1} \lambda_i^s(t+1)\, p_{j \rightarrow i}^s(t) \left( 1 - \bar{p}_i^s \right) \prod_{m \in \partial i \backslash j} \left( 1 - \alpha_{mi} \cdot p_{m \rightarrow i}^s(t) \right) \\
&- \sum_{s \in S} \sum_{t=0}^{T-1} \lambda_j^s(t+1)\, p_{i \rightarrow j}^s(t) \left( 1 - \bar{p}_j^s \right) \prod_{m \in \partial j \backslash i} \left( 1 - \alpha_{mj} \cdot p_{m \rightarrow j}^s(t) \right) \\
&- \sum_{s \in S} \sum_{t=0}^{T-1} \sum_{k \in \partial i \backslash j} \lambda_{i \rightarrow k}^s(t+1)\, p_{j \rightarrow i}^s(t) \left( 1 - \bar{p}_i^s \right) \prod_{m \in \partial i \backslash \{j,k\}} \left( 1 - \alpha_{mi} \cdot p_{m \rightarrow i}^s(t) \right) \\
&- \sum_{s \in S} \sum_{t=0}^{T-1} \sum_{k \in \partial j \backslash i} \lambda_{j \rightarrow k}^s(t+1)\, p_{i \rightarrow j}^s(t) \left( 1 - \bar{p}_j^s \right) \prod_{m \in \partial j \backslash \{i,k\}} \left( 1 - \alpha_{mj} \cdot p_{m \rightarrow j}^s(t) \right),
\end{aligned} \tag{B5}
$$

which can further be simplified for $\alpha_{ij} \neq 0$:

$$
\frac{\partial \mathcal{L}}{\partial \alpha_{ij}} = -\frac{1}{\alpha_{ij}} \sum_{s \in S} \sum_{t=0}^{T-1} \left( \lambda_{i \rightarrow j}^s(t)\, p_{i \rightarrow j}^s(t) + \lambda_{j \rightarrow i}^s(t)\, p_{j \rightarrow i}^s(t) \right). \tag{B6}
$$

The above can be directly used in the forward–backward propagation procedure described in the main text.

## 2. Simple Graphs

In the case of simple graphs, where $\alpha_{ij} = \alpha \,\forall_{(i,j)\in E}$, the objective is the same as for the original SLICER, but the constraints simplify to:

$$
\begin{aligned}
\mathcal{C} = &\sum_{s\in S}\sum_{t=0}^{T-1}\sum_{i\in V}\lambda_i^s(t+1)\left(p_i^s(t+1)-1+\left(1-\bar{p}_i^s\right)\prod_{k\in\partial i}\left(1-\alpha\cdot p_{k\to i}^s(t)\right)\right) \\
&+\sum_{s\in S}\sum_{t=0}^{T-1}\sum_{(i,j)\in E}\lambda_{i\to j}^s(t+1)\left(p_{i\to j}^s(t+1)-1+\left(1-\bar{p}_i^s\right)\prod_{k\in\partial i\setminus j}\left(1-\alpha\cdot p_{k\to i}^s(t)\right)\right).
\end{aligned}
\tag{B7}
$$

As a result, we need to update the derivatives over messages:

$$
\begin{aligned}
\frac{\partial\mathcal{L}}{\partial p_{i\to j}^s(t)} = &\,\lambda_{i\to j}^s(t)-\lambda_j^s(t+1)\,\alpha\left(1-\bar{p}_j^s\right)\prod_{m\in\partial j\setminus i}\left(1-\alpha\cdot p_{m\to j}^s(t)\right) \\
&-\sum_{k\in\partial j\setminus i}\lambda_{j\to k}^s(t+1)\,\alpha\left(1-\bar{p}_j^s\right)\prod_{m\in\partial j\setminus\{i,k\}}\left(1-\alpha\cdot p_{m\to j}^s(t)\right),
\end{aligned}
\tag{B8}
$$

while the derivatives over marginals remain in the same form as before. Then we can rewrite the derivatives over parameters, which in this case reduce to only a single equation:

$$
\begin{aligned}
\frac{\partial\mathcal{L}}{\partial\alpha} = &-\sum_{s\in S}\sum_{t=0}^{T-1}\sum_{i\in V}\lambda_i^s(t+1)\sum_{j\in\partial i}p_{j\to i}^s(t)\left(1-\bar{p}_i^s\right)\prod_{m\in\partial i\setminus j}\left(1-\alpha\cdot p_{m\to i}^s(t)\right) \\
&-\sum_{s\in S}\sum_{t=0}^{T-1}\sum_{(i,k)\in E'}\lambda_{i\to k}^s(t+1)\sum_{j\in\partial i\setminus k}p_{j\to i}^s(t)\left(1-\bar{p}_i^s\right)\prod_{m\in\partial i\setminus\{j,k\}}\left(1-\alpha\cdot p_{m\to i}^s(t)\right),
\end{aligned}
\tag{B9}
$$

where by $E'$ we denote the set of edges such that $(i,j)$ and $(j,i)$ are counted as separate elements. The above equation can, assuming that $\alpha > 0$, be simplified to:

$$
\frac{\partial\mathcal{L}}{\partial\alpha} = -\frac{1}{\alpha}\sum_{s\in S}\sum_{t=0}^{T-1}\sum_{(i,j)\in E'}\lambda_{i\to j}^s(t)\cdot p_{i\to j}^s(t).
\tag{B10}
$$

Now we can repeat the whole procedure, but the update step is in only one dimension.

## 3. Missing Times

In this setting part of the activation times are known precisely, while others are known to be in certain time intervals. Based on that we modify the objective in a following way.

$$
\begin{aligned}
\mathcal{O} = &\sum_{s\in S}\sum_{i\in O}\sum_{\tau_i^s}m^{\tau_i^s}\log\left(p_i^s(\tau_i^s)\cdot\mathbb{1}_{(\tau_i^s\leq T)}-p_i^s(\tau_i^s-1)\cdot\mathbb{1}_{(\tau_i^s>0)}+\mathbb{1}_{(\tau_i^s=*)}\right) \\
&+\sum_{s\in S}\sum_{i\in O}\sum_{\hat{\tau}_i^s}m^{\hat{\tau}_i^s}\log\left(p_i^s(\hat{\tau}_i^s)-p_i^s(\hat{\tau}_i^s-\Delta\hat{\tau}_i^s)\right).
\end{aligned}
\tag{B11}
$$

The only change is that if the node was activated during an unobserved time interval, the marginal probability of being infected precisely at $\tau_i$ is replaced by a marginal probability of being infected inside the interval $[\hat{\tau}_i-\Delta\hat{\tau}_i,\hat{\tau}_i]$. The constraints are the same as in the original SLICER. The above modification affects only the computation of the

Lagrange derivative over marginal probabilities.

$$
\frac{\partial \mathcal{L}}{\partial p_i^s(t)} = \lambda_i^s(t) + \sum_{\tau_i^s} \frac{m^{\tau_i^s} \cdot \mathbb{1}_{(t=\tau_i^s)} \cdot \mathbb{1}_{(\tau_i^s \leq T)}}{p_i^s(\tau_i^s) - p_i^s(\tau_i^s - 1) \cdot \mathbb{1}_{(\tau_i^s > 0)}}
$$

$$
+ \sum_{\tau_i^s} \frac{m^{\tau_i^s} \cdot \mathbb{1}_{(t=\tau_i^s - 1)} \cdot \mathbb{1}_{(\tau_i^s > 0)}}{p_i^s(\tau_i^s - 1) - p_i^s(\tau_i^s) \cdot \mathbb{1}_{(\tau_i^s \leq T)} - \mathbb{1}_{(\tau_i^s = *)}} \tag{B12}
$$

$$
+ \sum_{\hat{\tau}_i^s} \frac{m^{\hat{\tau}_i^s} \cdot \mathbb{1}_{(t=\hat{\tau}_i^s)}}{p_i^s(\hat{\tau}_i^s) - p_i^s(\hat{\tau}_i^s - \Delta\hat{\tau}_i^s)} + \sum_{\hat{\tau}_i^s} \frac{m^{\hat{\tau}_i^s} \cdot \mathbb{1}_{(t=\hat{\tau}_i^s - \Delta\hat{\tau}_i^s)}}{p_i^s(\hat{\tau}_i^s - \Delta\hat{\tau}_i^s) - p_i^s(\hat{\tau}_i^s)}.
$$

This requires updating the values of Lagrange multipliers $\lambda_i^s(t)$. The rest of the procedure remains the same.

## 4. Noisy Timestamps

Taking into account the noisy activation times results with an objective function of the following form:

$$
\mathcal{O} = \sum_{s \in S} \sum_{i \in O} \sum_{\tau_i^s} m^{\tau_i^s} \log\left( \sum_{k=-K}^{K} \pi_k \left[ p_i^s(\tau_i^s - k) \cdot \mathbb{1}_{(0 \leq \tau_i^s - k \leq T)} - p_i^s(\tau_i^s - 1 - k) \cdot \mathbb{1}_{(0 < \tau_i^s - k \leq *)} + \mathbb{1}_{(\tau_i^s - k = *)} \right] \right). \tag{B13}
$$

Updating the objective function affects the learning procedure. However, the only element of the optimisation scheme, which require changes is the computation of $\lambda_i^s(t)$.

$$
\frac{\partial \mathcal{L}}{\partial p_i^s(t)} = \lambda_i^s(t) + \sum_{\tau_i^s} \frac{m^{\tau_i^s} \cdot \sum_{k=-K}^{K} \pi_k \left( \mathbb{1}_{(t=\tau_i^s - k)} \cdot \mathbb{1}_{(0 \leq \tau_i^s - k \leq T)} - \mathbb{1}_{(t=\tau_i^s - 1 - k)} \cdot \mathbb{1}_{(0 < \tau_i^s - k \leq *)} \right)}{\sum_{k=-K}^{K} \pi_k \left( p_i^s(\tau_i^s - k) \cdot \mathbb{1}_{(0 \leq \tau_i^s - k \leq T)} - p_i^s(\tau_i^s - 1 - k) \cdot \mathbb{1}_{(0 < \tau_i^s - k \leq *)} + \mathbb{1}_{(\tau_i^s - k = *)} \right)}. \tag{B14}
$$

All the further computations remain the same as before. It should be noted that the described change in the learning procedure does also affect its complexity. In the worst case scenario where $K = T$ the complexity of the algorithm becomes $\max[O(|E||S|T), O(N|S|T^2)]$ in comparison with the complexity $O(|E||S|T)$ of the original algorithm.

## Appendix C: Efficient Implementation

Direct implementation of Eq. (2) and (3) suffers from evaluating the products over neighbors, in particular, this is repeated for every message rooted in a given node. As a result the complexity of the algorithm depends on the network's degree distribution, instead of the number of edges. Networks with broad degree distributions, which are common in real-world applications, are specifically affected by this problem. It is possible to overcome this by rewriting the DMP equations in the following form:

$$
p_i^c(t) = 1 - \left(1 - \bar{p}_i^c\right) \prod_{k \in \partial i} \left(1 - \alpha_{ki} \cdot p_{k \to i}^c(t-1)\right), \tag{C1}
$$

$$
p_{j \to i}^c(t) = 1 - \frac{1 - p_j^c(t)}{1 - \alpha_{ij} \cdot p_{i \to j}^c(t-1)}, \tag{C2}
$$

where the marginals equations are the same, but the messages are computed using the pre-computed marginals. Now the complexity is independent of the degree distribution and truly linear in the number of edges, assuming that $\alpha_{ij} \cdot p_{i \to j}^c(t-1) \neq 1$. The latter is, however, a reasonable assumption, which is unlikely to be broken, in which case the original equation can still be used. Similar problem arises from Eq. (B4), which after equating left hand side to zero leads to:

$$
\lambda_{i \to j}^s(t) = \lambda_j^s(t+1)\, \alpha_{ij} \left(1 - \bar{p}_j^s\right) \prod_{m \in \partial j \setminus i} \left(1 - \alpha_{mj} \cdot p_{m \to j}^s(t)\right)
$$

$$
+ \sum_{k \in \partial j \setminus i} \lambda_{j \to k}^s(t+1)\, \alpha_{ij} \left(1 - \bar{p}_j^s\right) \prod_{m \in \partial j \setminus \{i,k\}} \left(1 - \alpha_{mj} \cdot p_{m \to j}^s(t)\right). \tag{C3}
$$

Although the equation have a different form, they are dual to the DMP one. Similarly to the latter, direct implementation leads to degree distribution dependence. Here instead of a product over neighbors we have a sum over neighbors and similarly as before, this can be prevented. To do so, let us introduce

$$\widehat{\lambda}_j^s(t) = \sum_{k \in \partial j} \lambda_{j \to k}^s(t+1)\left(1 - \bar{p}_j^s\right) \prod_{m \in \partial j \setminus k} \left(1 - \alpha_{mj} \cdot p_{m \to j}^s(t)\right), \tag{C4}$$

which can be further simplified as

$$\widehat{\lambda}_j^s(t) = \sum_{k \in \partial j} \lambda_{j \to k}^s(t+1)\left(1 - p_{j \to k}^s(t+1)\right). \tag{C5}$$

Now we can write:

$$\begin{aligned} \lambda_{i \to j}^s(t) &= \lambda_j^s(t+1)\,\alpha_{ij}\left(1 - p_{j \to i}^s(t+1)\right) \\ &+ \alpha_{ij}\,\frac{\widehat{\lambda}_j^s(t) - \lambda_{j \to i}^s(t+1)\left(1 - p_{j \to i}^s(t+1)\right)}{1 - \alpha_{ij} \cdot p_{i \to j}^s(t)}. \end{aligned} \tag{C6}$$

This way, knowing that $\lambda_{i \to j}^s(T) = 0\;\forall_{i,j}$, we reduced the complexity to being linear in the number of edges, same way as we did for messages. Additionally, in the directed case where $(i,j) \in E$, but $(j,i) \notin E$, one needs to rewrite Eq. (C6) into:

$$\lambda_{i \to j}^s(t) = \alpha_{ij}\,\frac{\widehat{\lambda}_j^s(t) + \lambda_j^s(t+1)\left(1 - p_j^s(t+1)\right)}{1 - \alpha_{ij} \cdot p_{i \to j}^s(t)}. \tag{C7}$$

The above equations, used in all our experiments, are implemented in Julia [60] and available at [59].

## Appendix D: Structure learning without any prior knowledge on the topology

Here, we present results on structure learning in the presence of unobserved nodes. We reproduce the experiment presented in Fig. 4, but in the most challenging setting where no information on the network topology is available, and the super-set of edges that we initially consider corresponds to a fully-connected graph. The results are given in Fig. 10, and are similar for all four synthetic network types. The edge reconstruction is always better than random guess, but we observe reduction of reconstruction quality with the increase of the number of unobserved nodes. This is expected due to the appearance of degeneracy where several models are consistent with the observation data. As a result SLICER produces a graph, which is compatible with the observed data and as such can be used as an effective model for prediction [8], even if it is not the *true* graph. Moreover, have in mind that our method recovers the network simultaneously with learning model parameters.

[1] Mark Newman. *Networks*. Oxford university press, 2018.

[2] Alain Barrat, Marc Barthelemy, and Alessandro Vespignani. *Dynamical processes on complex networks*. Cambridge university press, 2008.

[3] Praneeth Netrapalli and Sujay Sanghavi. Learning the graph of epidemic cascades. *ACM SIGMETRICS Performance Evaluation Review*, 40(1):211–222, 2012.

[4] Manuel Gomez-Rodriguez, Jure Leskovec, and Andreas Krause. Inferring networks of diffusion and influence. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 5(4):1–37, 2012.

[5] Bruno Abrahao, Flavio Chierichetti, Robert Kleinberg, and Alessandro Panconesi. Trace complexity of network inference. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 491–499, 2013.

[6] Jean Pouget-Abadie and Thibaut Horel. Inferring graphs from cascades: A sparse recovery framework. In *International Conference on Machine Learning*, pages 977–986. PMLR, 2015.

[7] Andrey Y Lokhov. Reconstructing parameters of spreading models from partial observations. *Advances in Neural Information Processing Systems*, 29, 2016.

[8] Mateusz Wilinski and Andrey Y Lokhov. Prediction-centric learning of independent cascade dynamics from partial observations. In *International Conference on Machine Learning*, pages 11182–11192. PMLR, 2021.

[9] Isaak Neri and Désiré Bollé. The cavity approach to parallel dynamics of ising spins on a graph. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(08):P08009, 2009.
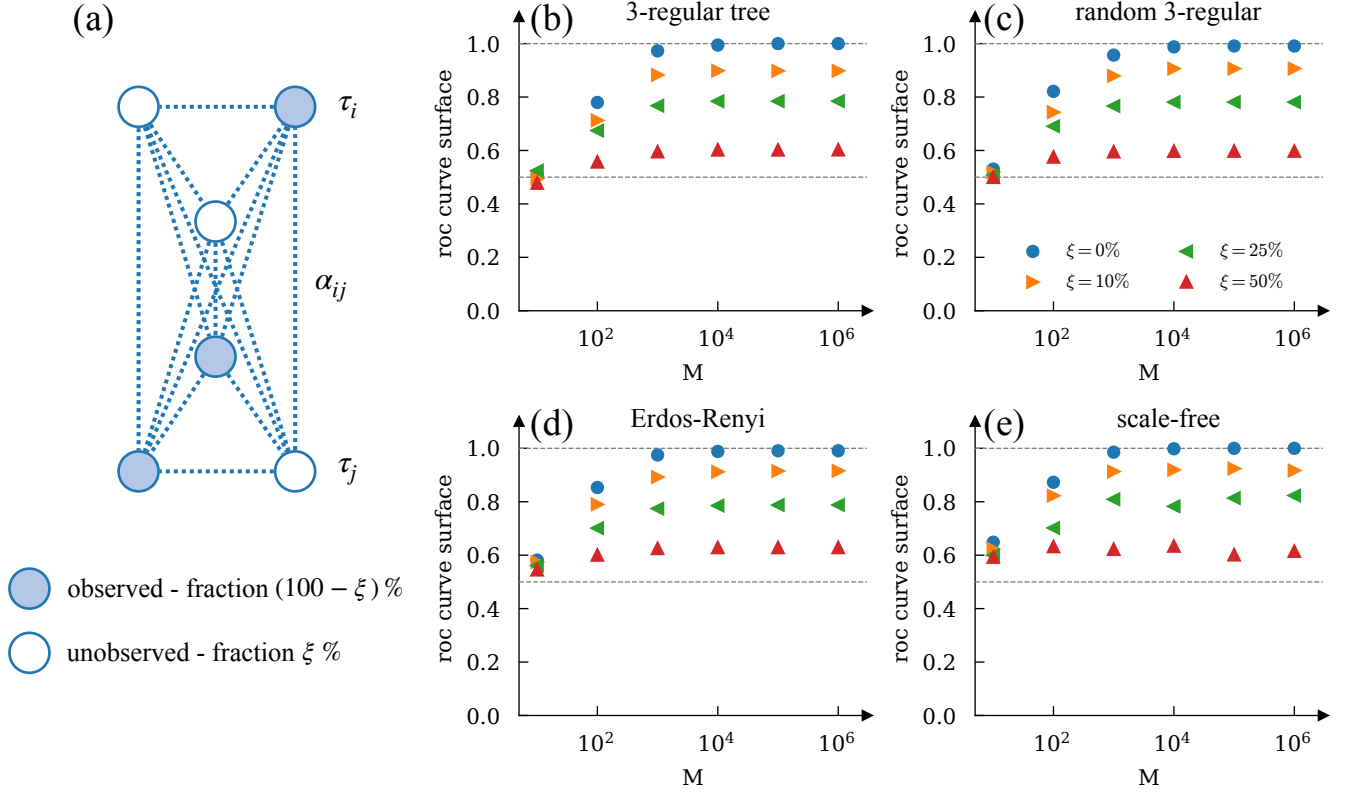
FIG. 10: Structure learning starting with a fully-connected graph. Average ROC curve surface, as a function of the number of available cascades for different network types in the case where a fraction $\xi$ of nodes is unobserved and there is no prior knowledge about the topology. Each point is averaged over 5 different networks and 5 different sets of parameters $\alpha_{ij}^*$ (sampled from a uniform distribution in $[0,1]$). All networks contain $N = 100$ nodes and all but the tree have average degree equal to $\langle k \rangle = 3$. Unobserved nodes were picked at random. All cascades had length $T = 5$.

[10] Brian Karrer and Mark EJ Newman. Message passing approach for general epidemic models. *Physical Review E*, 82(1):016101, 2010.

[11] Yashodhan Kanoria and Andrea Montanari. Majority dynamics on trees and the dynamic cavity method. 2011.

[12] Erik Aurell and Hamed Mahmoudi. Three lemmas on dynamic cavity method. *Communications in Theoretical Physics*, 56(1):157, 2011.

[13] Erik Aurell and Hamed Mahmoudi. Dynamic mean-field and cavity methods for diluted ising systems. *Physical Review E*, 85(3):031119, 2012.

[14] Fabrizio Altarelli, Alfredo Braunstein, Luca Dall'Asta, and Riccardo Zecchina. Large deviations of cascade processes on graphs. *Physical Review E*, 87(6):062115, 2013.

[15] Fabrizio Altarelli, Alfredo Braunstein, Luca Dall'Asta, and Riccardo Zecchina. Optimizing spread dynamics on graphs by message passing. *Journal of Statistical Mechanics: Theory and Experiment*, 2013(09):P09011, 2013.

[16] Andrey Y Lokhov, Marc Mézard, Hiroki Ohta, and Lenka Zdeborová. Inferring the origin of an epidemic with a dynamic message-passing algorithm. *Physical Review E*, 90(1):012801, 2014.

[17] Andrey Y Lokhov, Marc Mézard, and Lenka Zdeborová. Dynamic message-passing equations for models with unidirectional dynamics. *Physical Review E*, 91(1):012811, 2015.

[18] Gino Del Ferraro and Erik Aurell. Dynamic message-passing approach for kinetic spin models with reversible dynamics. *Physical Review E*, 92(1):010102, 2015.

[19] Munik Shrestha, Samuel V Scarpino, and Cristopher Moore. Message-passing approach for recurrent-state epidemic models on networks. *Physical Review E*, 92(2):022821, 2015.

[20] Andrey Y Lokhov and David Saad. Optimal deployment of resources for maximizing impact in spreading processes. *Proceedings of the National Academy of Sciences*, 114(39):E8138–E8146, 2017.

[21] Andrey Y Lokhov and David Saad. Scalable influence estimation without sampling. *arXiv preprint arXiv:1912.12749*, 2019.

[22] Thomas Barthel. The matrix product approximation for the dynamic cavity method. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(1):013217, 2020.

[23] Hanlin Sun, David Saad, and Andrey Y Lokhov. Competition, collaboration, and optimization in multiple interacting spreading processes. *Physical Review X*, 11(1):011048, 2021.

[24] Bo Li and David Saad. Impact of presymptomatic transmission on epidemic spreading in contact networks: A dynamic message-passing analysis. *Physical Review E*, 103(5):052303, 2021.

[25] Erik Aurell, David Machado Perez, and Roberto Mulet. A closure for the master equation starting from the dynamic cavity method. *Journal of Physics A: Mathematical and Theoretical*, 56(17):17LT02, 2023.

[26] Stefano Crotti and Alfredo Braunstein. Matrix product belief propagation for reweighted stochastic dynamics over graphs. *Proceedings of the National Academy of Sciences*, 120(47):e2307935120, 2023.

[27] Freya Behrens, Barbora Hudcová, and Lenka Zdeborová. Backtracking dynamical cavity method. *Phys. Rev. X*, 13:031021, Aug 2023.

[28] Davide Ghio, Antoine LM Aragon, Indaco Biazzo, and Lenka Zdeborová. Bayes-optimal inference for spreading processes on random networks. *arXiv preprint arXiv:2303.17704*, 2023.

[29] Jonathan S Yedidia, William T Freeman, Yair Weiss, et al. Understanding belief propagation and its generalizations. *Exploring artificial intelligence in the new millennium*, 8(236-239):0018–9448, 2003.

[30] Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.

[31] Jacob Goldenberg, Barak Libai, and Eitan Muller. Talk of the network: A complex systems look at the underlying process of word-of-mouth. *Marketing letters*, 12(3):211–223, 2001.

[32] David Kempe, Jon Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 137–146, 2003.

[33] William Ogilvy Kermack and Anderson G McKendrick. A contribution to the mathematical theory of epidemics. *Proceedings of the royal society of london. Series A, Containing papers of a mathematical and physical character*, 115(772):700–721, 1927.

[34] Matt J Keeling and Ken TD Eames. Networks and epidemic models. *Journal of the royal society interface*, 2(4):295–307, 2005.

[35] Hao Huang, Keqi Han, Beicheng Xu, and Ting Gan. Reconstructing diffusion networks from incomplete data. In *IJCAI*, volume 2022, pages 3085–3091, 2022.

[36] Daryl J Daley and Joe Gani. *Epidemic modelling: an introduction*. Number 15. Cambridge University Press, 2001.

[37] Emre Sefer and Carl Kingsford. Convex risk minimization to infer networks from probabilistic diffusion data at multiple scales. In *2015 IEEE 31st International Conference on Data Engineering*, pages 663–674. IEEE, 2015.

[38] Ritabrata Dutta, Antonietta Mira, and Jukka-Pekka Onnela. Bayesian inference of spreading processes on networks. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 474(2215):20180129, 2018.

[39] Kareem Amin, Hoda Heidari, and Michael Kearns. Learning from contagion (without timestamps). In *International Conference on Machine Learning*, pages 1845–1853. PMLR, 2014.

[40] Adisak Supeesun and Jittat Fakcharoenphol. Learning network structures from contagion. *Information Processing Letters*, 121:11–16, 2017.

[41] Alfredo Braunstein, Alessandro Ingrosso, and Anna Paola Muntoni. Network reconstruction from infection cascades. *Journal of the Royal Society Interface*, 16(151):20180844, 2019.

[42] Keqi Han, Yuan Tian, Yunjia Zhang, Ling Han, Hao Huang, and Yunjun Gao. Statistical estimation of diffusion network topologies. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pages 625–636. IEEE, 2020.

[43] Jessica Hoffmann and Constantine Caramanis. Learning graphs from noisy epidemic cascades. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(2):1–34, 2019.

[44] William Trouleau, Jalal Etesami, Matthias Grossglauser, Negar Kiyavash, and Patrick Thiran. Learning hawkes processes under synchronization noise. In *International Conference on Machine Learning*, pages 6325–6334. PMLR, 2019.

[45] Wei Chen, Chi Wang, and Yajun Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1029–1038, 2010.

[46] Michael Shapiro and Edgar Delgado-Eckert. Finding the probability of infection in an sir network is np-hard. *Mathematical biosciences*, 240(2):77–84, 2012.

[47] Wei Chen, Yajun Wang, and Siyu Yang. Efficient influence maximization in social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 199–208, 2009.

[48] Nan Du, Le Song, Manuel Gomez Rodriguez, and Hongyuan Zha. Scalable influence estimation in continuous-time diffusion networks. *Advances in neural information processing systems*, 26, 2013.

[49] Edith Cohen, Daniel Delling, Thomas Pajor, and Renato F Werneck. Sketch-based influence maximization and computation: Scaling up with guarantees. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pages 629–638, 2014.

[50] Brendan Lucier, Joel Oren, and Yaron Singer. Influence at scale: Distributed computation of complex contagion in networks. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 735–744, 2015.

[51] Hung T Nguyen, Tri P Nguyen, Tam N Vu, and Thang N Dinh. Outward influence and cascade size estimation in billion-scale networks. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(1):1–30, 2017.

[52] Béla Bollobás. Random graphs. In *Modern graph theory*, pages 215–252. Springer, 1998.

[53] Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.

[54] Tom Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006.

[55] Mark EJ Newman. The structure and function of complex networks. *SIAM review*, 45(2):167–256, 2003.

[56] Wayne W Zachary. An information flow model for conflict and fission in small groups. *Journal of anthropological research*, 33(4):452–473, 1977.

[57] Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.

[58] Amanda L Traud, Peter J Mucha, and Mason A Porter. Social structure of Facebook networks. *Phys. A*, 391(16):4165–4180, Aug 2012.

[59] Mateusz Wilinski and Andrey Y Lokhov. https://github.com/mateuszwilinski/dynamic-message-passing/, 2021.

[60] Stefan Karpinski Jeff Bezanson. The julia programming language.