



HAL
open science

Compensating for population sampling in simulations of epidemic spread on temporal contact networks

Mathieu Génois, Christian L. Vestergaard, Ciro Cattuto, Alain Barrat

► **To cite this version:**

Mathieu Génois, Christian L. Vestergaard, Ciro Cattuto, Alain Barrat. Compensating for population sampling in simulations of epidemic spread on temporal contact networks. *Nature Communications*, 2015, 6, pp.8860. hal-01131855

HAL Id: hal-01131855

<https://hal.science/hal-01131855>

Submitted on 16 Mar 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Compensating for population sampling in simulations of epidemic spread on temporal contact networks

Mathieu Génois,¹ Christian L. Vestergaard,¹ Ciro Cattuto,² and Alain Barrat^{1,2}

¹*Aix Marseille Université, Université de Toulon,
CNRS, CPT, UMR 7332, 13288 Marseille, France*

²*Data Science Laboratory, ISI Foundation, Torino, Italy*

(Dated: March 13, 2015)

Data describing human interactions often suffer from incomplete sampling of the underlying population. As a consequence, the study of contagion processes using data-driven models can lead to a severe underestimation of the epidemic risk. Here we present a systematic method to correct this bias and obtain an accurate estimation of the risk in the context of epidemic models informed by high-resolution time-resolved contact data. We consider several such data sets collected in various contexts and perform controlled resampling experiments. We show that the statistical information contained in the resampled data allows us to build surrogate versions of the unknown contacts and that simulations of epidemic processes using these surrogate data sets yield good estimates of the outcome of simulations performed using the complete data set. We discuss limitations and potential improvements of our method.

Human interactions play an important role in determining the potential transmission routes of infectious diseases and other contagion phenomena [1]. Their measure and characterisation thus represent an invaluable contribution notably to the study of transmissible diseases [2]. In this context, the use of surveys and diaries in which volunteer participants record their encounters [3–8] have provided crucial insight, despite the memory biases inherent in self-reporting procedures [4, 9, 10]. Moreover, new approaches have emerged to measure contact patterns between individuals, using wearable sensors that can detect the proximity of other similar devices [11–20]. Data gathering efforts have produced data sets describing the contact patterns between individuals in various contexts in the form of temporal networks [15, 17, 21–24]: nodes represent individuals and, at each time step, a link is drawn between pairs of individuals who are in contact [25]. Such data can inform models of epidemic spreading phenomena to evaluate epidemic risks and mitigation strategies [15, 22, 26–31].

However, most data sets suffer from population sampling: despite efforts to maximise participation, for instance through scientific engagement of participants [24, 32], not all individuals accept to participate. Hence, the collected data only contains information on contacts occurring among a fraction of the population under study.

Population sampling is known to affect the properties of static networks [33, 34]: Various statistical properties of a sampled network may differ from those of the complete system under scrutiny [35], and several works have focused on inferring network statistics from the knowledge of incomplete, sampled network data [36–39]. Both structural and temporal properties of time-varying networks are as well affected by sampling effects.

In addition, a crucial though poorly studied consequence of population sampling is that simulations of dynamical processes in data-driven models can be affected. For instance, in simulations of epidemic spreading, ex-

cluded nodes are by definition unreachable and thus equivalent to immunised nodes. Due to herd vaccination effects, the outcome of simulations of epidemic models on sampled networks is thus underestimated. How to estimate the outcome of dynamical processes on contact networks using incomplete data remains an open question.

Here we tackle this issue for incompletely sampled data describing networks of human face-to-face interactions. We do not aim at inferring the true sequence of missing contacts but at estimating the outcome of simulations of models of epidemic spread in the whole population. To this effect, we resample available data sets by excluding at random a fraction of the individuals (nodes of the contact network), measure how resampling affects relevant network statistics and show that some crucial properties are stable under resampling. We exploit this stability to present a systematic method to construct surrogate contact sequences for the excluded nodes, using only information available in the resampled data. We show that the outcome of simulations performed on the reconstructed data sets, obtained by the union of the resampled and surrogate contacts, reproduces results obtained on the complete data set, while using only the resampled data severely underestimates the epidemic risk. We show the efficiency of our procedure for data collected in three widely different contexts: a conference, a high school and a workplace.

RESULTS

Data & Methodology

We consider data sets describing contacts between individuals, collected by the SocioPatterns collaboration [17] in three different settings: an office building (“InVS”) [40], a high school (“Thiers13”) [24] and a scientific con-

ference (“SFHH”) [21, 22]. These data correspond to the close face-to-face proximity of individuals equipped with wearable sensors, at a temporal resolution of 20 seconds [17, 18]. Table I summarises the characteristics of each data set. The contact data are represented by temporal networks, in which nodes represent the participating individuals and a link between two nodes i and j at time t indicates that the two corresponding persons were in contact at that time. Moreover, the InVS and Thiers13 populations were structured in departments and classes, respectively (see Methods).

To understand how sampling affects the properties of the temporal networks, the outcome of simulations of dynamical processes, and, ultimately, how to compensate for these effects, we consider as ground truth the available data¹ and perform population resampling experiments by removing a fraction f of the nodes uniformly at random. We explore how several characteristics of the temporal networks depend on f and how the outcome of numerical simulations of epidemic spread is biased by such resampling. We then present a method for constructing surrogate data using only information contained in the resampled data and show that it allows us to obtain a good estimate of the outcome of processes simulated on the complete data set.

Specifically, we consider the Susceptible-Infectious-Recovered (SIR) and the Susceptible-Infectious-Susceptible (SIS) models of epidemic propagation. In these models, a susceptible (S) node becomes infectious (I) at rate β when in contact with an infectious node. Infectious nodes recover spontaneously at rate μ . In the SIR model, nodes then enter an immune *recovered* (R) state, while in the SIS model, nodes become susceptible again and can thus be reinfected. The quantities of interest are for the SIR model the distribution of epidemic sizes, defined as the final fraction of recovered nodes, and for the SIS model the average fraction of infectious nodes i_∞ in the stationary state. We also calculate for the SIR model the fraction of epidemics that infect more than 20% of the population and the average size of these epidemics. For the SIS model, we determine the *epidemic threshold* β_c for different values of μ : it corresponds to the value of β that separates an epidemic-free state ($i_\infty = 0$) for $\beta < \beta_c$ from an endemic state ($i_\infty > 0$) for $\beta > \beta_c$, and is thus an important indicator of the epidemic risk. We refer to the Methods section for further details on models and simulations.

To validate our method, we compare for each data set the outcomes of simulations performed (i) on the whole data set, (ii) on resampled data sets with a varying fraction of nodes removed, f , and (iii) on reconstructed data sets built to compensate for sampling effects.

Resampled contact networks

Sampling affects the various properties of contact networks in different ways. The number of neighbours of a node decreases as the fraction f of removed nodes increases (Supplementary Fig. S1) since removing nodes also removes links to these nodes. On the contrary, the density of the resampled aggregated network, defined as the number of links divided by the total number of possible links between the nodes, remains stable (Supplementary Fig. S2).

We moreover show in Supplementary Figs. S2–S4 that the group structure of a population, as summarised by contact matrices, is well preserved under resampling. Contact matrices give a measure of the interaction between groups (here classes or departments). Here we consider the *link density* matrix, where the element (i, j) is given by the number of links between individuals of groups i and j in the aggregated contact network, normalised by the total number of possible links between these two groups². Table II and Supplementary Fig. S2 show that the similarities between the original and sampled matrices are high for all data sets (see Supplementary Figs. S3–S4 for the contact matrices themselves).

Finally, the temporal statistics of the contact network are not affected by population sampling, as already noted in [18] for other data sets: the distributions of contact durations, of inter-contact durations, of number of contacts per link and of cumulated contact durations (i.e., of the link weights in the aggregated network) do not change when the network is sampled (Supplementary Fig. S1).

Despite the robustness of these properties under sampling, the outcome of simulations of epidemic spread is strongly affected by the resampling (Figs. 1 and 2). In particular, for large values of f , the probability of large outbreaks in the SIR model vanishes (Fig. 1). As mentioned above, the reason is that the removed nodes act as if they were immunised: sampling hinders the propagation by removing transmission routes between the remaining nodes. As a result, the prevalence and the final size of the outbreaks are systematically underestimated by simulations of the SIR model on the resampled network with respect to simulations on the whole data set (Fig. 1), and the epidemic threshold of the SIS model is overestimated (Fig. 2): resampling leads to a systematic underestimation of the epidemic risk.

Estimation of epidemic sizes through simulations on reconstructed temporal networks

We now develop a method able to compensate for the bias due to population sampling in the simulations of

¹ Note that the full data sets are also samples of all the contacts that occurred in the populations, as the participation rate was lower than 100% in each case.

² If n_i denotes the number of individuals in group i , the number of possible links is equal to $n_i n_j / 2$ for $i \neq j$ and to $n_i(n_i - 1) / 2$ for $i = j$.

epidemic spread. To this aim, we use only information measured in the resampled network that has been shown above to be robust under resampling, namely the density of the aggregated contact network, the contact matrices of link density and the distributions of number of contacts per link and of contact and inter-contact durations. We assume that the number of missing nodes in each group is known, as well as the timelines of scheduled activity (nights and weekends, during which no contact occurs).

For each data set, considered as ground truth, we create resampled data sets by removing at random a fraction f of the N nodes. We measure on each resampled data set the contact network density, the contact matrix of link densities between groups (for the structured population cases) as well as the distributions mentioned above. We then construct stochastic, surrogate versions of the missing part of the network, as described in detail in the Methods section: We create for each missing node a surrogate instance of its links, in a way to keep the density of the network and the density of links between groups of nodes fixed (equal to the value measured in the resampled data); we then create a synthetic timeline of contacts on each surrogate link, extracted at random from the measured distributions. The statistical properties of the resulting reconstructed (surrogate) networks, obtained by the union of the resampled data and of the surrogate links, are similar to the ones of the original data (Table II and Supplementary Figs. S5-S10). In particular, the method conserves the distribution of link weights (W), the group structure (S), and the temporal (T) characteristics of individual links; we thus refer to it in the following as the **WST** method. We emphasise that our aim is not to infer the true missing contacts, and we do not compare the detailed structures of the surrogate and original contact networks. Instead, we simply aim at obtaining a surrogate version of these contacts, such that the resulting simulations of epidemic spread yield outcomes close to the ones obtained when the whole data set is used.

Figure 1 shows that the distribution of final epidemic sizes for the SIR model on surrogate networks is much closer to those obtained by simulating on the whole data set than for the simulations performed on the resampled networks. Instead of a severe underestimation of epidemic sizes, we obtain a slight overestimation. We quantify this result in Fig. 3 where we display the fraction of outbreaks that reach at least 20% of the population and the average epidemic size for these outbreaks, for simulations performed on either resampled or surrogate networks. In the case of resampled data, we rapidly lose information about the size and even the existence of large outbreaks, whereas this crucial information is well recovered when using surrogate networks, even when a large fraction of nodes are removed. We show in the Supplementary Information that similar results are obtained for different values of the spreading parameters. Moreover, as shown in Fig. 2 and Supplementary Figs. S11-S12, the phase diagram obtained for the SIS model when us-

ing reconstructed networks is much closer to the original than for resampled networks, in particular concerning the value of the epidemic threshold. Overall, simulations using the reconstructed network yield a much better estimation of the epidemic risk than simulations using resampled network data, for both SIS and SIR models.

When the fraction f of nodes excluded by the resampling procedure becomes large, the properties of the resampled data may start to differ substantially from those of the whole data set (Figs. S1 & S2). As a result, the distributions of epidemic sizes of SIR simulations deviate from those obtained on the whole data set (Fig. 4), even if the epidemic risk evaluation is still better than for simulations on the resampled networks. Most importantly however, the information remaining in the resampled data at large f can be insufficient to construct surrogate contacts. This happens in particular if an entire class or department is absent from the resampled data or if all the resampled nodes of a class/department are disconnected (see Supplementary Information for details). We investigate this limitation in the bottom plots of Fig. 3 by showing the failure rate, i.e., the fraction of cases in which we are not able to construct surrogate networks from the resampled data. The failure rate increases gradually with f for the InVS data since the groups (departments) are of different sizes. For the Thiers13 data, all classes are of similar sizes so that the failure rate reaches abruptly a large value at a given value of f . For the SFHH data, we can always construct surrogate networks as the population is not structured.

Reconstruction methods using partial information

A natural question concerns how much information is needed to build the surrogate networks. The **WST** method uses the contact matrix of link densities and the distributions of contact durations, inter-contact durations and numbers of contacts per link. We investigate here three alternative procedures that use less information, still only computed from the resampled data:

- **W**: we construct surrogate networks that conserve the overall link density and the distribution of weights found in the resampled data. The group structure (S) and temporal characteristics (T) of contacts are not conserved.
- **WS**: in addition to the **W** construction, we here use the contact matrix of link densities to conserve the group structure (S) of the network.
- **WT**: we construct the surrogate links as in the **W** method, without taking into account the link density contact matrix (S), but with contact timelines on each surrogate link that conserve the temporal characteristics (T) of the individual links.

Details of these methods are given in the Supplementary Information.

The results of SIR simulations performed on the resulting surrogate networks are shown in Fig. 3. All three methods using partial information lead to a larger overestimation of the epidemic risk than the **WST** method. The **W** method consistently gives the worst results, as infections spread easier on random homogeneous graphs than on structured graphs [41, 42]. Taking into account the population structure yields slightly better results (**WS** case), while using realistic contact sequences (i.e., taking into account the heterogeneity of contact numbers and durations and the burstiness of contacts) has an even stronger effect (**WT**). Overall, surrogate networks that respect all these constraints (**WST**) yield the best results.

DISCUSSION

The understanding of epidemic spreading phenomena have been vastly improved thanks to the use of data-driven models. In the case of contact data, population sampling represents however a serious issue: individuals absent from a data set have the same role as immunised individuals when epidemic processes are simulated. Feeding sampled data into data-driven models can therefore lead to severe underestimations of the epidemic risk and might even a priori affect the evaluation of mitigation strategies.

Here we have put forward a systematic method to compensate for such underestimation in simulations performed using sampled contact data and to obtain a good estimate of the epidemic risk in the entire population, as measured by the distribution of outbreak sizes and by the epidemic threshold. To this aim, we have shown how it is possible, starting from a data set describing the contacts of only a fraction of the population of interest, to construct a surrogate data set that uses only accessible information, i.e., quantities measured in the sampled data: the contact network density, the densities of links between groups in a structured population, and the distributions of numbers and durations of contacts and of inter-contact durations. Simulations of epidemic spreading on such surrogate data yield results similar to those obtained on the complete data set, albeit with a small overestimation of the risk (which, from a public health point of view, is much better than a large underestimation).

Strikingly, the method presented here yields good results even when a substantial part of the population is excluded, in particular in estimating the probability of large outbreaks. It has however two limitations when the fraction f of excluded individuals is too large. First, the construction of the surrogate contacts relies on the stability of a set of quantities with respect to resampling, but the measured quantities start to deviate from the original ones at large f . The shape of the distribution of epidemic sizes may then deviate substantially from the original one. Second, large values of f might even render

the construction of the surrogate data impossible due to the loss of information on whole categories of nodes.

Another important point concerns the issue of how much information should be included when constructing the surrogate data. It is linked to the general issue of how much information is needed to get an accurate picture of spreading processes on temporal networks [22, 27, 29, 43, 44]. On the one hand, we have shown that using less information to build the surrogate data yields worse results. On the other hand, the surrogate contacts do not take into account any correlations between structure and activity in the temporal contact network, which are known to influence spreading processes [21, 42, 44–47]. This might be the cause of the systematic overestimation of the epidemic sizes in simulations using reconstructed data (Figs. 1 and 3). From this point of view, our method could be improved by integrating into the surrogate data complex correlation patterns measured in the sampled data. It might for instance be possible to use the temporal network decomposition technique put forward in [47, 48] on the sampled data, in order to extract mesostructures such as temporally-localized mixing patterns. The surrogate contacts could then be built in a way to preserve such patterns.

Finally, we have considered uniform sampling of nodes, corresponding to data collection in a population with a participation rate smaller than 100%. It would be interesting to consider as well other types of data losses, due for instance to a partial coverage of the premises of interest by the measuring infrastructure [17]. Moreover, the population under study is (usually) not isolated from the external world, and it would be important to devise ways to include contacts with outsiders in the data and simulations, for instance by using other data sources such as surveys.

ACKNOWLEDGEMENTS

The present work is partially supported by the French ANR project HarMS-flu (ANR-12-MONU-0018) to M.G. and A.B., by the EU FET project Multiplex 317532 to A.B., C.C. and C.L.V., by the A*MIDEX project (ANR-11-IDEX-0001-02) funded by the "Investissements d'Avenir" French Government program, managed by the French National Research Agency (ANR) to A.B., by the Lagrange Project of the ISI Foundation funded by the CRT Foundation to C.C., and by the Q-ARACNE project funded by the Fondazione Compagnia di San Paolo to C.C. We warmly thank André Panisson for providing the algorithm used in the numerical simulations of the SIS process.

AUTHOR CONTRIBUTIONS

A.B. and C.C. designed and supervised the study. M.G., C.L.V., A.B., and C.C. collected and post-

processed the data. M.G., C.L.V. and A.B. analyzed the data, carried out computer simulations and prepared the figures. M.G., C.L.V., A.B., and C.C. wrote the manuscript.

METHODS

Data

Two of the data sets used in our study – the high school (Thiers13) and the office building (InVS) – are structured in groups corresponding respectively to classes and departments. For the conference data (SFHH), we instead do not have metadata on the participants, and it has been shown in [22] that the aggregated network structure was homogeneous.

The high school is structured in 9 classes, which form three subgroups of three classes corresponding to their specialisation in Mathematics-Physics (MP, MP*1, MP*2 with respectively 31, 29 and 38 students), Physics (PC, PC*, PSI with respectively 44, 39 and 34 students), or Biology (2BIO1, 2BIO2, 2BIO3 with respectively 37, 35 and 39 students).

The office building is structured in 5 departments: DISQ (Scientific Direction, 15 persons), DMCT (Department of Chronic Diseases and Traumatism, 26 persons), DSE (Department of Health and Environment, 34 persons), SRH (Human Resources, 13 persons) and SFLE (Logistics, 4 persons).

SIR simulations

In the SIR model, all nodes are initially susceptible, except one infectious node, the seed of the epidemic. A susceptible node in contact with an infectious node becomes infectious at rate β . Infectious nodes enter the recovered state at rate μ .

Spreading simulations are performed using the temporal networks of contacts (original, resampled or reconstructed). We run each simulation until no infectious individual remains (nodes are thus either still susceptible or have been infected and then recovered). If needed, the sequence of contacts is repeated in the simulation [22]. We consider values of β and μ yielding a non-negligible epidemic risk, i.e., such that a rather large fraction of simulations lead to a final size larger than 20% of the population (see Fig. 1): $\beta = 4 \times 10^{-4} s^{-1}$, $\mu = 4 \times 10^{-7} s^{-1}$ (InVS) or $4 \times 10^{-6} s^{-1}$ (SFHH and Thiers13). Other parameter values are explored in the Supplementary Information. For each set of parameters, the distribution of epidemic sizes is obtained by performing 1,000 simulations.

SIS phase diagram

In the SIS model infectious nodes become susceptible again at rate μ . Depending on the values of β and μ , the outbreak can either die out or reach a stationary state. For each value of μ , there is a critical threshold value for β where the system goes continuously from a stationary state with no infectious nodes, to a stationary state with a finite fraction of infectious nodes. In order to determine this threshold, we follow [49] and write the Markov chain equation for the probability $p_i^{(t)}$ for a node to be infectious at time t :

$$p_i^{(t)} = 1 - \left(1 - (1 - \mu)p_i^{(t-1)}\right) \prod_j \left(1 - \beta A_{ji}^{(t-1)} p_j^{(t-1)}\right). \quad (1)$$

Here $A^{(t)}$ is the adjacency matrix of the network at the time t ($A_{ij}^{(t)} = 1$ if there is a link between nodes i and j at time t , and $A_{ij}^{(t)} = 0$ otherwise). For each data set of finite length T , we assume periodic boundary conditions with $A_{ij}^{(t+T)} = A_{ij}^{(t)}$. This equation corresponds to the extension of the individual-based mean field approach³ [50] to SIS processes on temporal networks. The equation is iterated until a stationary state is reached, and we then compute the average fraction of infectious nodes as

$$i_\infty = \sum_{i,t} \frac{p_i^{(t)}}{TN}, \quad (2)$$

where T is the duration of the data set and N the number of nodes.

Reconstruction algorithm

We consider a population \mathcal{P} of N individuals (the nodes of the contact network), potentially organised in groups. We assume that all the contacts taking place among a subpopulation $\tilde{\mathcal{P}}$ of these individuals, of size $\tilde{N} = (1 - f)N$, are known. This constitutes our resampled data from which we need to construct a surrogate set of contacts concerning the remaining $n = N - \tilde{N} = fN$ individuals for which no contact information is available: these contacts can occur among these individuals and between them and the members of $\tilde{\mathcal{P}}$. We assume that we know the group to which each member of $\mathcal{P} \setminus \tilde{\mathcal{P}}$ belongs, as well as the overall activity timeline, i.e. the intervals during which contacts take place, separated by nights and weekends. We construct the surrogate contacts as follows:

1. we measure in the sampled data:

³ It does not take into account the dynamical correlations created by the spreading process between the states of neighboring nodes.

- the density ρ of links in the time-aggregated network;
 - a row-normalised contact matrix C , in which the element C_{AB} gives the probability for a node in group A to have a link to a node of group B ;
 - the list $\{\tau_c\}$ of contact durations;
 - the list $\{\tau_{ic}\}$ of inter-contact durations;
 - the list $\{p\}$ of numbers of contacts per link;
2. we compute the number e of additional links needed to keep the network density constant when we add the n excluded nodes.
 3. we construct each link according to the following procedure:
 - a node i is randomly chosen from the set $\mathcal{P} \setminus \tilde{\mathcal{P}}$ of excluded nodes;
 - knowing the group A that i belongs to, we extract at random a target group B with probability given by C_{AB} ;
 - we draw a target node j at random from B (if $B = A$, we take care that $i \neq j$) such that i and j are not linked;
 - from $\{p\}$, we draw the number of contact events p taking place over the link ij ;
 - the starting time t_0 of the first ij contact is drawn uniformly in the first activity interval of the global activity timeline;
 - we then sequentially draw the duration of a contact from $\{\tau_c\}$ and the inter-contact duration until the next contact from $\{\tau_{ic}\}$, and repeat until we have built p contacts;
 - finally, we insert breaks defined by the activity timeline.

When the fraction f of excluded nodes is large, the method may be unable to construct the surrogate links (see Supplementary Information).

-
- [1] Barrat, A., Barthélemy, M. & Vespignani, A. *Dynamical processes on complex networks* (Cambridge University Press (Cambridge), 2008).
 - [2] Read, J. M., Edmunds, W. J., Riley, S., Lessler, J. & Cummings, D. A. T. Close encounters of the infectious kind: methods to measure social mixing behaviour. *Epidemiology & Infection* **140**, 2117–2130 (2012).
 - [3] Edmunds, W. J., O’callaghan, C. J. & Nokes, D. J. Who mixes with whom? a method to determine the contact patterns of adults that may lead to the spread of airborne infections. *Proceedings of the Royal Society of London. Series B: Biological Sciences* **264**, 949–957 (1997).
 - [4] Read, J., Eames, K. & Edmunds, W. Dynamic social networks and the implications for the spread of infectious disease. *J R Soc Interface* **5**, 1001–7 (2008).
 - [5] Mossong, J. *et al.* Social contacts and mixing patterns relevant to the spread of infectious diseases. *PLoS Med* **5**, e74 (2008).
 - [6] Mikolajczyk, R., Akmatov, M., Rastin, S. & Kretzschmar, M. Social contacts of school children and the transmission of respiratory-spread pathogens. *Epidemiology & Infection* **136**, 813–822 (2008).
 - [7] Danon, L., House, T., Read, J. & Keeling, M. Social encounter networks: collective properties and disease transmission. *J. R. Soc. Interface* **9**, 2826–2833 (2012).
 - [8] Danon, L., Read, J. M., House, T. A., Vernon, M. C. & Keeling, M. J. Social encounter networks: characterizing great britain. *Proceedings of the Royal Society B: Biological Sciences* **280** (2013).
 - [9] Smieszek, T., Burri, E. U., Scherzinger, R. & Scholz, R. W. Collecting close-contact social mixing data with contact diaries: reporting errors and biases. *Epidemiology & Infection* **140**, 744–752 (2012).
 - [10] Smieszek, T. *et al.* How should social mixing be measured: comparing web-based survey and sensor-based methods. *BMC Infectious Diseases* **14**, 136 (2014).
 - [11] Hui, P. *et al.* Pocket switched networks and human mobility in conference environments. In *WDTN ’05: Proc. 2005 ACM SIGCOMM workshop on Delay-tolerant networking* (ACM, New York, NY, USA, 2005).
 - [12] O’Neill, E. *et al.* Instrumenting the city: Developing methods for observing and understanding the digital cityscape. In *UbiComp*, 315–332 (2006).
 - [13] Eagle, N., Pentland, A. S. & Lazer, D. Inferring friendship network structure by using mobile phone data. *Proceedings of the National Academy of Sciences* **106**, 15274–15278 (2009).
 - [14] Vu, L., Nahrstedt, K., Retika, S. & Gupta, I. Joint bluetooth/wifi scanning framework for characterizing and leveraging people movement in university campus. In *Proceedings of the 13th ACM International Conference on Modeling, Analysis, and Simulation of Wireless and Mobile Systems, MSWIM ’10*, 257–265 (ACM, New York, NY, USA, 2010).
 - [15] Salathé, M. *et al.* A high-resolution human contact network for infectious disease transmission. *Proceedings of the National Academy of Sciences* **107**, 22020–22025 (2010).
 - [16] Hashemian, M., Stanley, K. & Osgood, N. Flunet: Automated tracking of contacts during flu season. In *Proceedings of the 6th International workshop on Wireless Network Measurements*, 557–562 (2010).
 - [17] <http://www.sociopatterns.org>.
 - [18] Cattuto, C. *et al.* Dynamics of person-to-person interactions from distributed RFID sensor networks. *PLoS ONE* **5**, e11596 (2010).
 - [19] Hornbeck, T. *et al.* Using sensor networks to study the effect of peripatetic healthcare workers on the spread of hospital-associated infections. *Journal of Infectious Diseases* (2012).

- [20] Stopczynski, A. *et al.* Measuring large-scale social networks with high resolution. *PLoS ONE* **9**, e95978 (2014).
- [21] Isella, L. *et al.* What's in a crowd? Analysis of face-to-face behavioral networks. *Journal of Theoretical Biology* **271**, 166–180 (2011).
- [22] Stehlé, J. *et al.* Simulation of an SEIR infectious disease model on the dynamic contact network of conference attendees. *BMC Medicine* **9**, 87 (2011).
- [23] Stehlé, J. *et al.* High-resolution measurements of face-to-face contact patterns in a primary school. *PLoS ONE* **6**, e23176 (2011).
- [24] Fournet, J. & Barrat, A. Contact patterns among high school students. *PLoS ONE* **9**, e107878 (2014).
- [25] Holme, P. & Saramki, J. Temporal networks. *Physics Reports* **519**, 97 – 125 (2012).
- [26] Lee, S., Rocha, L. E. C., Liljeros, F. & Holme, P. Exploiting temporal network structures of human interaction to effectively immunize populations. *PLoS ONE* **7**, e36439 (2012).
- [27] Machens, A. *et al.* An infectious disease model on empirical networks of human contact: bridging the gap between dynamic network data and contact matrices. *BMC Infectious Diseases* **13**, 185 (2013).
- [28] Starnini, M., Machens, A., Cattuto, C., Barrat, A. & Pastor-Satorras, R. Immunization strategies for epidemic processes in time-varying contact networks. *Journal of theoretical biology* **337**, 89–100 (2013).
- [29] Smieszek, T. & Salathé, M. A low-cost method to assess the epidemiological importance of individuals in controlling infectious disease outbreaks. *BMC MEDICINE* **11**, 35 (2013). See related commentary article here <http://www.biomedcentral.com/1741-7015/11/36>.
- [30] Masuda, N. & Holme, P. Predicting and controlling infectious disease epidemics using temporal networks. *F1000Prime Reports* **5** (2013).
- [31] Gemmetto, V., Barrat, A. & Cattuto, C. Mitigation of infectious disease at school: targeted class closure vs school closure. *BMC Infectious Diseases* **14**, 695 (2014).
- [32] Conlan, A. J. K. *et al.* Measuring social networks in british primary schools through scientific engagement. *Proceedings of the Royal Society B: Biological Sciences* **278**, 1467–1475 (2011).
- [33] Granovetter, M. Network sampling: Some first steps. *American Journal of Sociology* **81**, pp. 1287–1303 (1976).
- [34] Frank, O. Sampling and estimation in large social networks. *Social Networks* **1**, 91 – 101 (1978/1979).
- [35] Kossinets, G. Effects of missing data in social networks. *Social Networks* **28**, 247 – 268 (2006).
- [36] Viger, F., Barrat, A., Dall'Asta, L., Zhang, C.-H. & Kolaczyk, E. What is the real size of a sampled network? the case of the Internet. *Phys. Rev. E* **75**, 056111 (2007).
- [37] Bliss, C. A., Danforth, C. M. & Dodds, P. S. Estimation of global network statistics from incomplete data. *PLoS ONE* **9**, e108471 (2014).
- [38] Zhang, Y., Kolaczyk, E. D. & Spencer, B. D. Estimating Network Degree Distributions Under Sampling: An Inverse Problem, with Applications to Monitoring Social Media Networks. *ArXiv e-prints* (2013). arXiv:1305.4977.
- [39] Cimini, G., Squartini, T., Gabrielli, A. & Garlaschelli, D. Systemic risk analysis in reconstructed economic and financial networks. *ArXiv e-prints* (2014). arXiv:1411.7613.
- [40] Géniois, M. *et al.* Data on face-to-face contacts in an office building suggests a low-cost vaccination strategy based on community linkers. *ArXiv e-prints* (2014). arXiv:1409.7017.
- [41] Onnela, J.-P. *et al.* Structure and tie strengths in mobile communication networks. *Proceedings of the National Academy of Sciences* **104**, 7332–7336 (2007).
- [42] Karsai, M. *et al.* Small but slow world: How network topology and burstiness slow down spreading. *Phys. Rev. E* **83**, 025102 (2011).
- [43] Blower, S. & Go, M.-H. The importance of including dynamic social networks when modeling epidemics of airborne infections: does increasing complexity increase accuracy? *BMC Medicine* **9**, 88 (2011).
- [44] Pfitzner, R., Scholtes, I., Garas, A., Tessone, C. J. & Schweitzer, F. Betweenness preference: Quantifying correlations in the topological dynamics of temporal networks. *Physical Review Letters* **110**, 198701 (2013).
- [45] Gauvin, L., Panisson, A., Cattuto, C. & Barrat, A. Activity clocks: spreading dynamics on temporal networks of human contact. *Scientific reports* **3** (2013).
- [46] Scholtes, I. *et al.* Causality-driven slow-down and speed-up of diffusion in non-markovian temporal networks. *Nat. Comm* **5**, 5024 (2014).
- [47] Gauvin, L., Panisson, A., Barrat, A. & Cattuto, C. Revealing latent factors of temporal networks for mesoscale intervention in epidemic spread. *ArXiv e-prints* (2015). arXiv:1501.02758.
- [48] Gauvin, L., Panisson, A. & Cattuto, C. Detecting the community structure and activity patterns of temporal networks: a non-negative tensor factorization approach. *PLoS ONE* **9**, e86028 (2014).
- [49] Valdano, E., Ferreri, L., Poletto, C. & Colizza, V. Analytical computation of the epidemic threshold on temporal networks. *ArXiv e-prints* (2014). arXiv:1406.4815.
- [50] Pastor-Satorras, R., Castellano, C., Van Mieghem, P. & Vespignani, A. Epidemic processes in complex networks. *ArXiv e-prints* (2014). arXiv:1408.2701.

TABLES & FIGURES

Data set	Type	N	r	T	Dates
InVS	Office building	92	63 %	2 weeks	June 24th - July 5th 2013
Thiers13	High school	326	86 %	1 week	December 2nd - 7th 2013
SFHH	Conference	403	34 %	2 days	June 3rd - 4th 2009

TABLE I. **Data sets considered.** For each data set we specify the type of social situation, the number N of individuals whose contacts were measured, the corresponding participation rate r , the duration T and the dates of the data collection.

	f	InVS CML	Thiers13 CML
Resampled	10 %	0.996 [0.937, 0.999]	0.999 [0.998, 0.999]
	20 %	0.980 [0.889, 0.994]	0.996 [0.995, 0.997]
	40 %	0.925 [0.872, 0.983]	0.988 [0.983, 0.990]
Reconstructed	10 %	0.976 [0.846, 0.995]	0.998 [0.994, 0.999]
	20 %	0.942 [0.844, 0.984]	0.993 [0.985, 0.995]
	40 %	0.890 [0.652, 0.953]	0.977 [0.938, 0.987]

TABLE II. **Similarities between the original contact matrices and the contact matrices of the resampled networks (top) and of the reconstructed networks (bottom).** Median and 90% confidence interval for the cosine similarity between link density contact matrices (CML) for different values of f , the fraction of nodes removed from the original data. Values were obtained from 100 independent realisations of the resampling and reconstruction procedures.

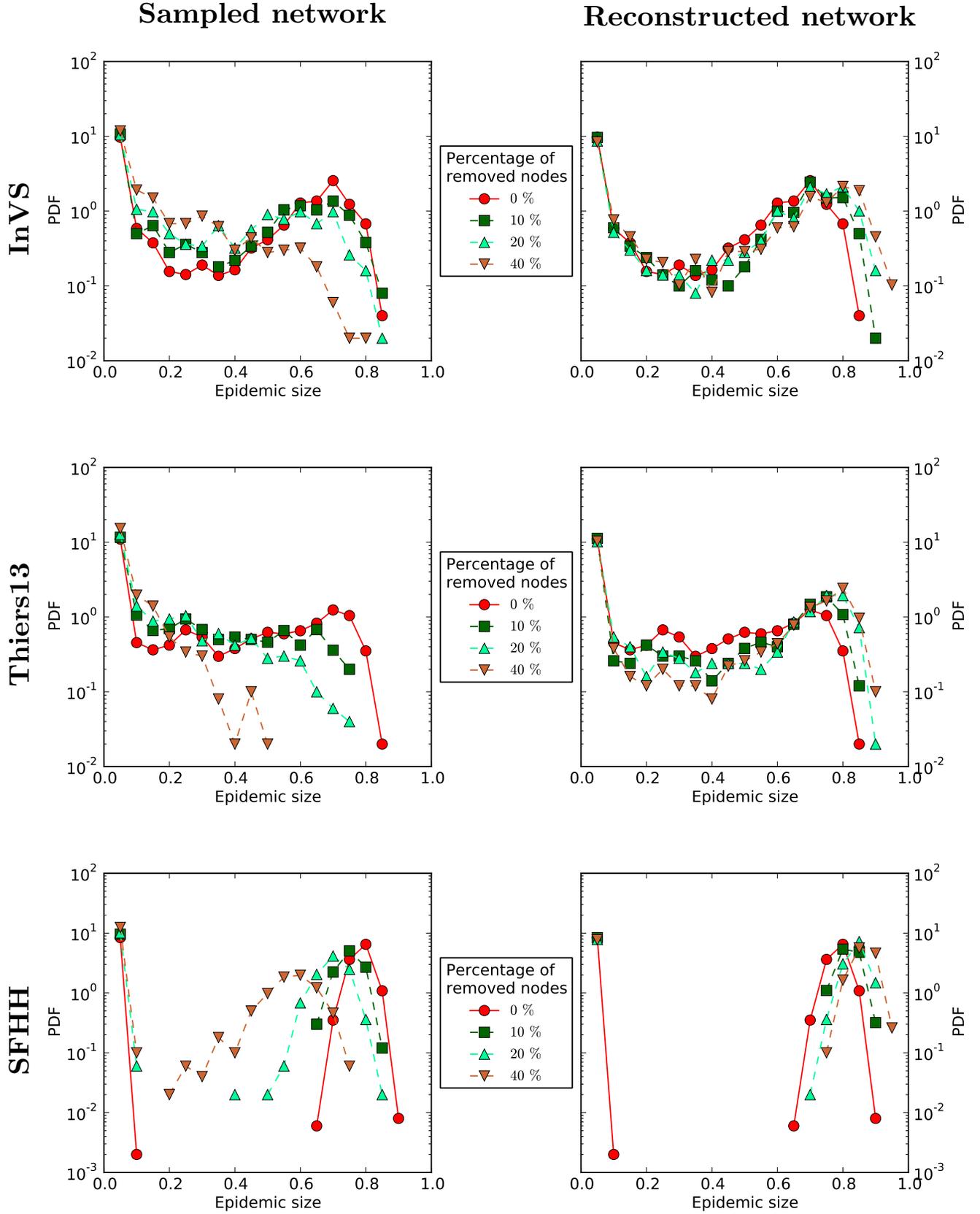


FIG. 1. **Comparison of the outcome of SIR epidemic simulations performed on resampled and reconstructed contact networks.** Distribution of epidemic sizes (fraction of recovered individuals) at the end of SIR processes simulated on top of either resampled (left column) or reconstructed (right) contact networks, for different values of the fraction f of nodes removed. The parameters of the SIR models are $\beta = 0.0004$ and $\beta/\mu = 1000$ (InVS) or $\beta/\mu = 100$ (Thiers13 and SFHH). The case $f = 0$ corresponds to simulations using the whole data set, i.e., the reference case. For each value of f , 1,000 independent simulations were performed.

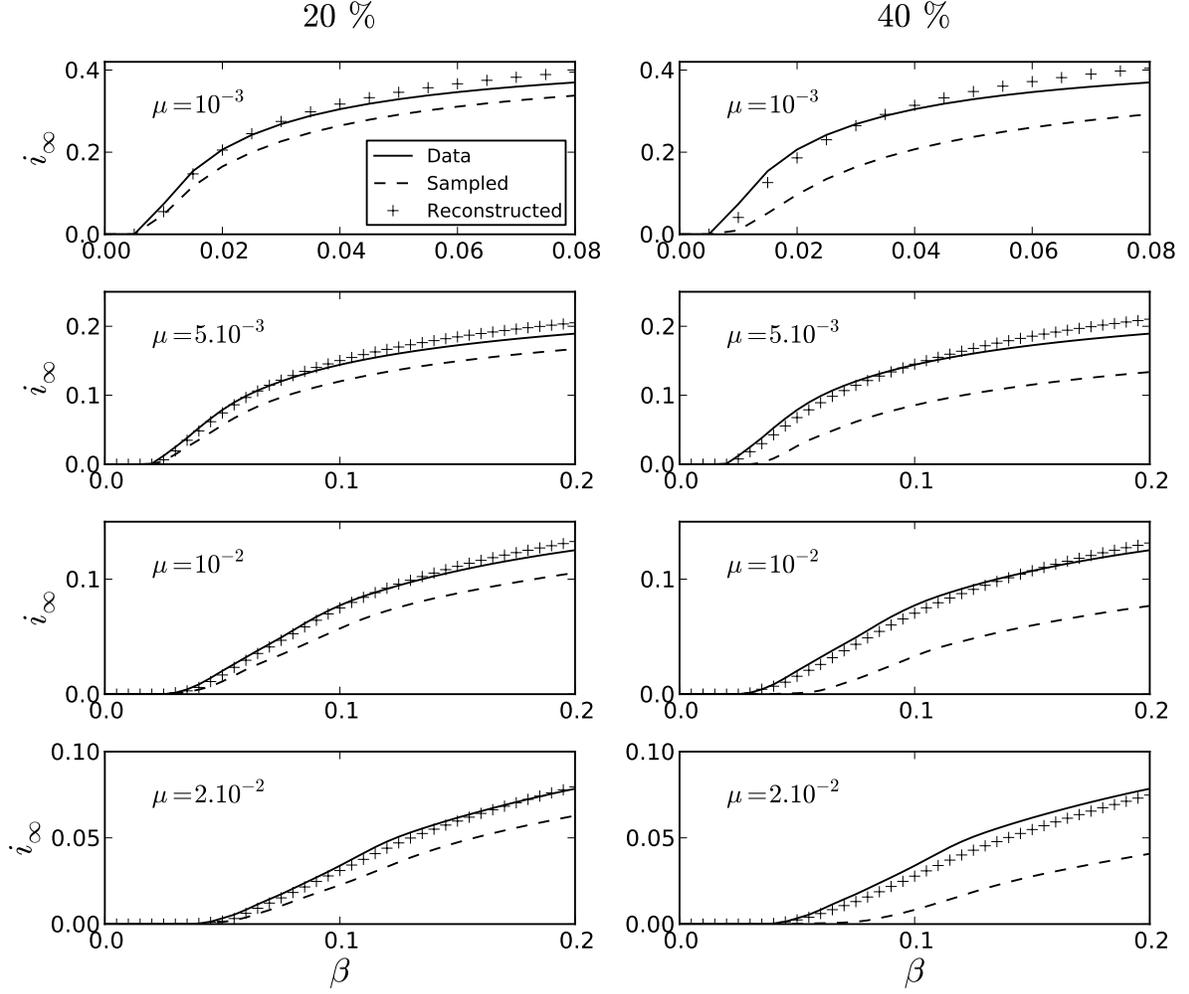


FIG. 2. **Phase diagram of the SIS model for the original, resampled and reconstructed contact networks (Thiers13 data set).** Each panel shows the stationary value i_∞ of the prevalence in the stationary state of the SIS model, computed as described in Methods (Eq. (2)), as a function of β , for several values of μ . Here we consider the example of the Thiers13 data set. The epidemic threshold corresponds to the transition between $i_\infty = 0$ and $i_\infty > 0$. Equations (1)-(2) are computed in each case using either the complete data set (continuous lines), resampled data (dashed lines) or reconstructed contact networks (pluses). The fraction of excluded nodes in the resampling is $f = 20\%$ for the left column and $f = 40\%$ for the right column.

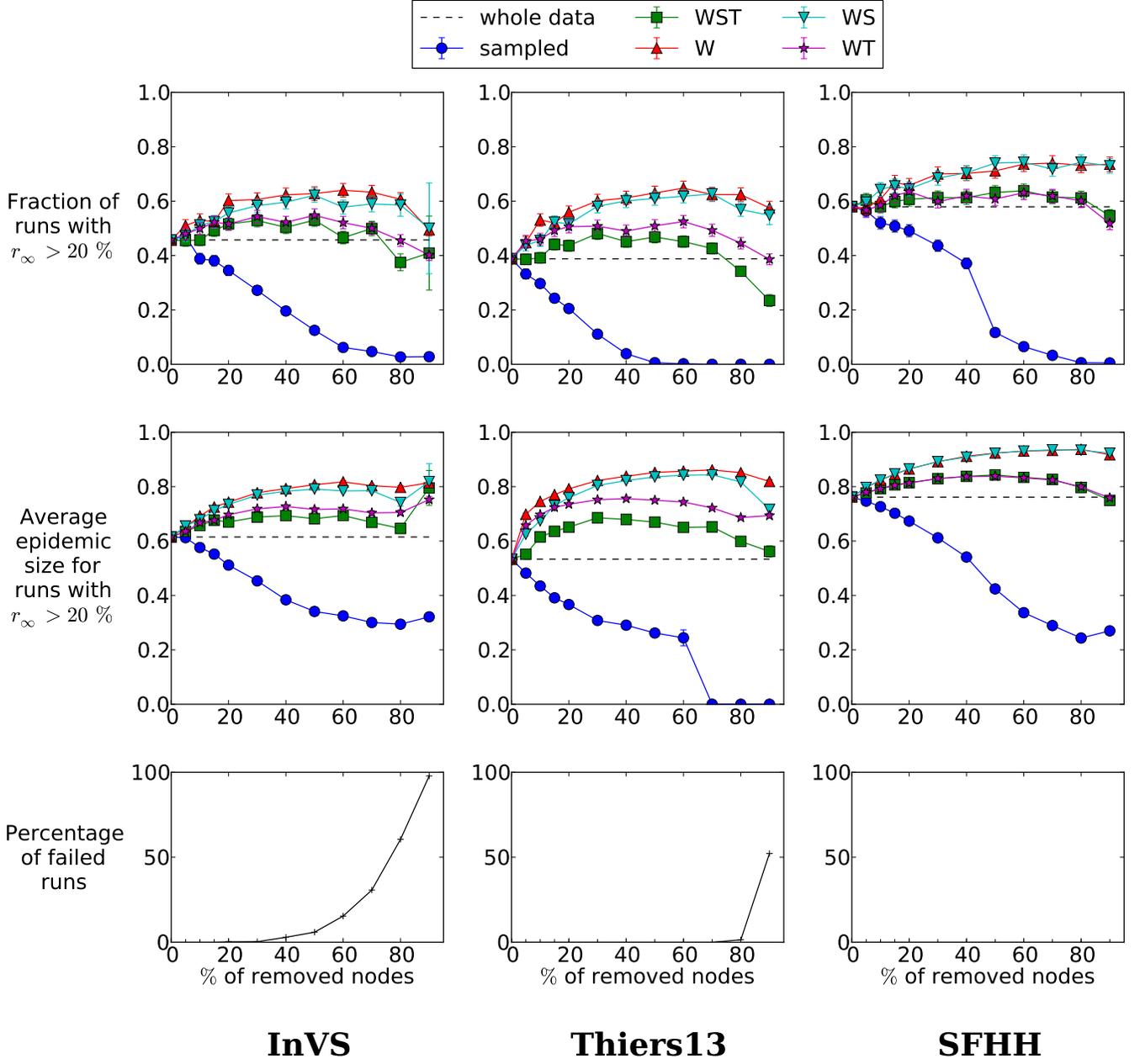


FIG. 3. **Outcome of SIR epidemic simulations performed on resampled contact networks and on networks reconstructed using different methods.** We compare in each case, and as a function of f , the fraction of outbreaks that lead to a final fraction of recovered individuals r_∞ larger than 20% of the population (top plots), and the average size of these large outbreaks (middle plots). The dashed lines give the corresponding values for simulations performed on the complete data sets. The different methods are (see text): reconstruction conserving only the link density and the distribution of weights of the resampled data (**W**); reconstruction preserving, in addition to the **W** method, the group structure of the resampled data (**WS**); reconstruction conserving link density, distribution of weights and distributions of contact times, of inter-contact times and of numbers of contacts per link measured in the resampled data (**WT**); full method conserving all these properties (**WST**). We also plot as a function of f the failure rate of the **WST** algorithm, i.e. the percentage of failed reconstructions (bottom plots); in the SFHH case reconstruction is always possible as the population is not structured into groups. The SIR parameters are $\beta = 0.0004$ and $\beta/\mu = 1000$ (InVS) or $\beta/\mu = 100$ (Thiers13 and SFHH) and each point is averaged over 1,000 independent simulations.

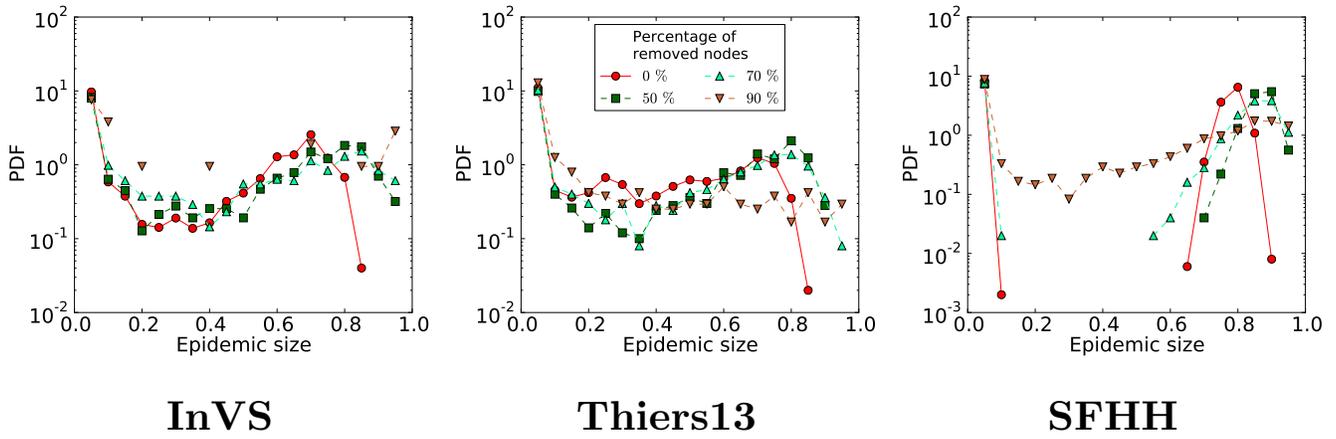


FIG. 4. **Outcome of simulations of SIR processes unfolding on reconstructed contact networks for large values of the fraction f of removed nodes.** Distributions of epidemic sizes for simulations of SIR processes on reconstructed networks and on the whole data set (case $f = 0$), similarly to Fig. 1, but for large values of the fraction f of removed nodes. Here $\beta = 0.0004$ and $\beta/\mu = 1000$ (InVS) or $\beta/\mu = 100$ (Thiers13 and SFHH) and 1,000 simulations were performed for each value of f . The distributions of epidemic sizes for simulations performed on resampled data sets are not shown since at these high values of f , almost no epidemics occur.

SUPPLEMENTARY INFORMATION

Effect of sampling on the temporal network of contacts

As described in the main text, we consider temporally resolved networks of contacts \mathcal{T} in a population \mathcal{P} of N individuals and we perform a resampling experiment by selecting a subpopulation $\tilde{\mathcal{P}}$ of these individuals, of size $\tilde{N} = (1 - f)N$. We assume that only the contacts occurring among the subpopulation $\tilde{\mathcal{P}}$ are known and we compare the properties of the corresponding resampled subnetwork $\tilde{\mathcal{T}}$ with those of the original network.

Figure S1 shows how population sampling affects several statistical properties of the contact networks. On the one hand, the degree distribution of the aggregated network of contacts systematically shifts towards smaller degree value. This is expected as each remaining node has in the resampled network a degree which is at most its degree in the original network, and is strictly smaller if some of its neighbours are not part of the resampled population. On the other hand, the statistical distributions of several quantities of interest are not affected by sampling: This is the case of the quantities attached either to single contacts or to single links, namely contact and inter-contact durations, number of contacts per link and link weights (the weight of a link is given by the total duration of the contacts between the two corresponding nodes).

Moreover, as shown in Fig. S2, the density of the aggregated network, i.e. the ratio between the number of links and the number of possible links, is on average conserved by the random resampling procedure. It varies however for different realisations of the resampling, and the corresponding variance increases with the fraction f of excluded nodes.

In the case of structured populations, Figures S3 & S4 show that the stability of the resampled network's density holds at the more detailed level of the contact matrices of link densities. In such matrices, the element (i, j) is given by the number of links between individuals of groups i and j , normalised by the total number of possible links between these two groups (if n_i denotes the number of individuals in group i , the number of possible links is equal to $n_i n_j / 2$ for $i \neq j$ and to $n_i(n_i - 1) / 2$ for $i = j$). These figures clearly illustrate how the diagonal and block-diagonal structures are preserved, and Fig. S2 gives a quantitative assessment of this stability by showing that the cosine similarity between contact matrices between the resampled and original aggregated contact networks remains high even for when a large fraction of the nodes are excluded.

Properties of the reconstructed contact networks

As described in the main text and in particular in the Methods section, we construct a surrogate set of con-

tacts concerning the fN individuals excluded by the resampling. We compare here the properties of the resulting contact networks (obtained by merging the resampled contact network $\tilde{\mathcal{T}}$ and the surrogate set of contacts) and of the original contact network. \mathcal{T} .

Figure S5 shows that the degree distribution, which is not constrained by the reconstruction procedure, deviates from the original distribution. On the other hand, the distributions of contact durations, inter-contact durations, number of contacts per link and link weights are preserved. Moreover, the link density contact matrices of the reconstructed networks (Fig. S6 & S7) share a high similarity with the original contact matrices, even for high fractions of nodes excluded (Fig. S10).

For completeness, we also compute the contact matrices in contact time density (CMT), in which each element (i, j) is given by the total time in contact between individuals of groups i and j , normalised by the total number of possible links between these two groups: it gives the average time spent in contact by two random individuals of groups i and j . Figures S8, S9 and S10 show that the structure of these matrices is well recovered by the reconstruction method, with high similarity with the original matrices.

Phase diagram of the SIS model for the conference and office building data sets

We observe for the workplace and the conference the same effect on the phase diagram of the SIS model as in the high school: sampling leads to a shift of the epidemic threshold to higher values and thus to an underestimation of the epidemic risk. The phase diagram and the epidemic threshold are estimated more accurately by using reconstructed networks, thus giving a better evaluation of the epidemic risk (Figs S11 & S12).

Sensitivity analysis

In the main text, we have considered values of the SIR model parameters leading to a non-negligible epidemic risk and a value of β corresponding to slow processes. We consider here several other values of the parameters, corresponding either to smaller epidemic risk (Fig. S13) or to faster processes (Figs. S14 - S16). In all cases, simulations performed on the resampled contact networks lead to a strong underestimation of the epidemic sizes, with distributions shifting to smaller values as f increases, while the use of reconstructed data sets leads to a better estimation and generally speaking a slight overestimation of the epidemic risk. The estimation of the distribution of epidemic sizes becomes worse when the processes become faster, which might be due to the fact that temporal correlations, which are present in the original data but not

in the surrogate contacts, play then a more important role in the outcome of the spreading processes.

Detailed alternative reconstruction methods

We give here details on the alternative reconstruction methods mentioned in the main text, which use less information than the **WST** method. In each case we consider the same setup as the complete method: a population \mathcal{P} of N individuals (the nodes of the contact network), potentially organised in groups, for which we know all the contacts taking place among a subpopulation $\tilde{\mathcal{P}}$ of size $\tilde{N} = (1 - f)N$. For the remaining $n = N - \tilde{N} = fN$ individuals, no contact information is available, but we know to which group they belong. We also have access to the overall activity timeline, *i.e.* to the successive intervals during which contacts can happen (daytimes), and are excluded (nights and weekends). The alternative reconstruction methods are the following:

W: We perform the reconstruction using only the network density and the distribution of link weights, both measured in the resampled network $\tilde{\mathcal{T}}$. The algorithm goes as follows:

1. we measure in the resampled data:
 - the density ρ of links in the time-aggregated network;
 - the list $\{w\}$ of link weights (the weight of a link is defined as the total contact time between the two linked nodes);
2. we compute the number of links e that must be added to keep the network density constant when we add the n excluded nodes;
3. we construct e links according to the following procedure:
 - a node i is randomly chosen from the set $\mathcal{P} \setminus \tilde{\mathcal{P}}$ of excluded nodes;
 - a node j is randomly chosen from the set $\mathcal{P} \setminus \{i\}$ of all other nodes;
 - from $\{w\}$, we draw the weight w of the link ij ;
 - we compute $n_{ij} = w/\Delta t$, where $\Delta t = 20s$ is the temporal resolution of the data set, and we randomly choose n_{ij} time windows of length Δt within the activity windows defined by the activity timeline as contact events between i and j .

WS: We perform the reconstruction using the network density, the distribution of link weights and the structure of the aggregated network given by the link density contact matrix, all measured in the resampled network $\tilde{\mathcal{T}}$. The algorithm goes as follows:

1. we measure in the resampled data:

- the density ρ of links in the time-aggregated network;
 - a *row-normalised* contact matrix C , in which the element C_{AB} gives the probability for a node in group A to have a link to a node of group B ;
 - the list $\{w\}$ of link weights;
2. we compute the number of links e that must be added to keep the network density constant when we add the n excluded nodes;
 3. we construct e links according to the following procedure:
 - a node i is randomly chosen from the set $\mathcal{P} \setminus \tilde{\mathcal{P}}$ of excluded nodes;
 - knowing the group A that i belongs to, we extract at random a target group B with probability given by C_{AB} ;
 - we draw a target node j at random from B (if $B = A$, we check that $j \neq i$);
 - from $\{w\}$, we draw the weight w of the link ij ;
 - as for the **W** method, we extract at random $w/\Delta t$ contact events of length $\Delta t = 20s$ within the activity timeline.

WT: We perform the reconstruction using the network density, the distribution of link weights and the temporal structure of the contacts given by the distributions of contact durations, inter-contact durations and number of contacts per link, all measured in the resampled network $\tilde{\mathcal{T}}$. The algorithm goes as follows:

1. we measure in the resampled data:
 - the density ρ of links in the time-aggregated network;
 - the list $\{\tau_c\}$ of contact durations;
 - the list $\{\tau_{ic}\}$ of inter-contact durations;
 - the list $\{p\}$ of numbers of contacts per link;
2. we compute the number of links e that must be added to keep the network density constant when we add the n excluded nodes;
3. we construct e links according to the following procedure:
 - a node i is randomly chosen from the set $\mathcal{P} \setminus \tilde{\mathcal{P}}$ of excluded nodes;
 - a node j is randomly chosen from the set $\mathcal{P} \setminus \{i\}$ of all other nodes;
 - from $\{p\}$, we draw the number of contact events p taking place over the link ij ;
 - the starting time t_0 of the first contact between i and j is drawn uniformly in the first activity interval of the global activity timeline;

- we then sequentially draw the duration of a contact from $\{\tau_c\}$ and the inter-contact duration until the next contact from $\{\tau_{ic}\}$, and repeat until we have built p contacts;
- finally, we insert breaks defined by the activity timeline.

Possible failure of the reconstruction method at large f

The construction of the surrogate version of the missing links uses as an input the group structure of the subgraph that remains after sampling, as given by the contact matrix of the link densities between the different groups of nodes that are present in $\tilde{\mathcal{P}}$. Depending on the characteristics of $\tilde{\mathcal{P}}$ and of the corresponding contacts, the construction method can fail in several cases:

- (i) if an entire group (class/department) of nodes in the population is absent from $\tilde{\mathcal{P}}$;
- (ii) if the remaining nodes of a specific group

(class/department) are all isolated in $\tilde{\mathcal{P}}$'s contact network;

- (iii) if, during the algorithm, a node of $\mathcal{P} \setminus \tilde{\mathcal{P}}$ is selected but cannot create any more links because it already has links to all nodes in the groups B such that $C_{AB} \neq 0$.

Cases (i) and (ii) correspond to a complete loss of information about the connectivity of a group (class/department) of the population, due to sampling. It is then impossible to reconstruct a sensible connectivity pattern for these nodes. Case (iii) is more subtle and occurs in situations of very low connectivity between groups. For instance, within the contact network of \mathcal{P} , a group A has links only with another specific group B , and both A and B are small; it is then possible that the nodes of $(\mathcal{P} \setminus \tilde{\mathcal{P}}) \cap A$ quickly exhaust the set of possible links to nodes of B during the reconstruction algorithm. If a node of $(\mathcal{P} \setminus \tilde{\mathcal{P}}) \cap A$ is again chosen to create a link, such a creation is not possible and the construction fails.

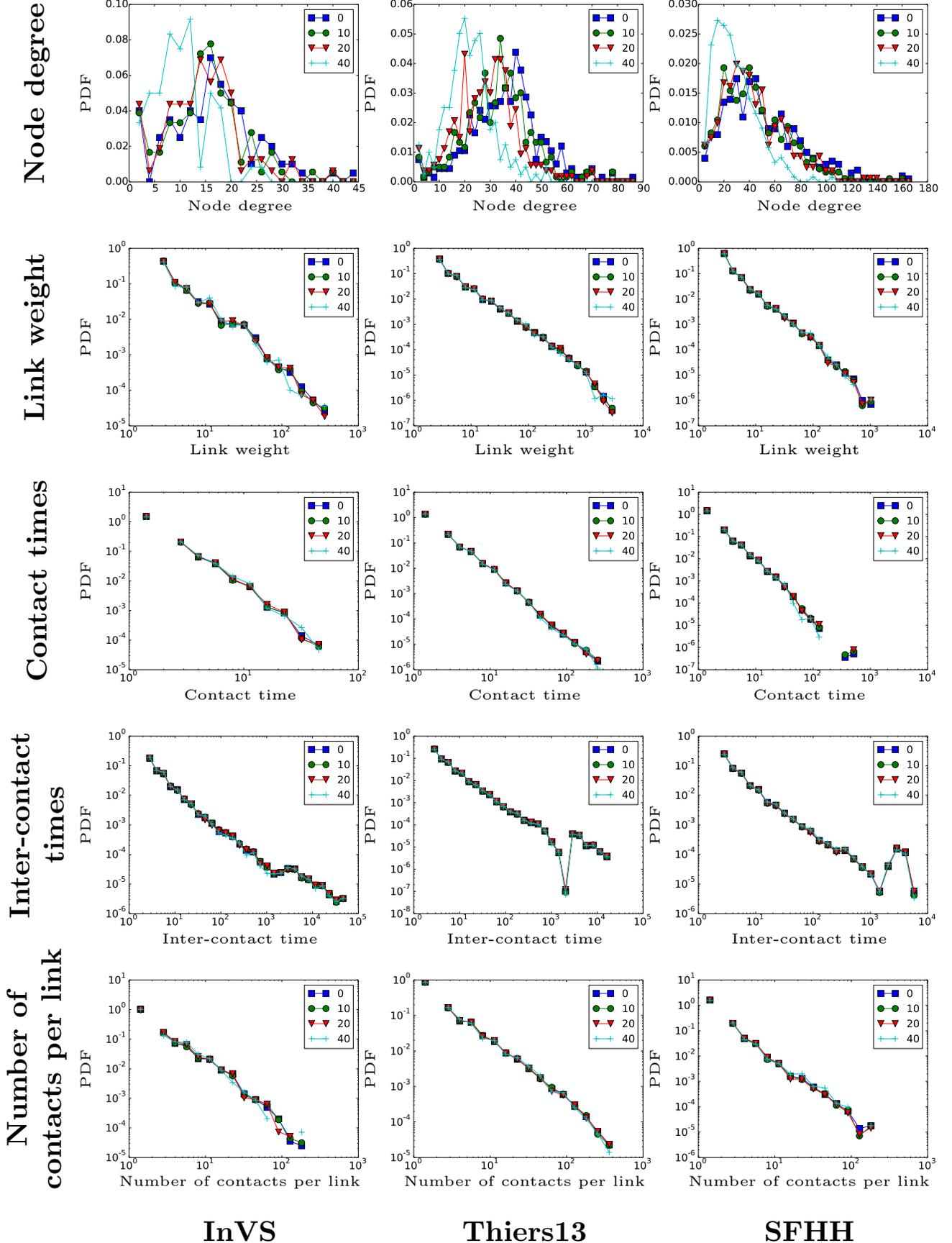


FIG. S1. **Effect of sampling on contact network properties.** Comparison of the distributions of structural (node degrees and link weights in the aggregated network of contacts) and temporal (contact durations, inter-contact times, number of contacts per link) properties of the contact networks, for different fractions f of removed nodes. For each value of f , the distributions are computed on a single realisation of the resampling.

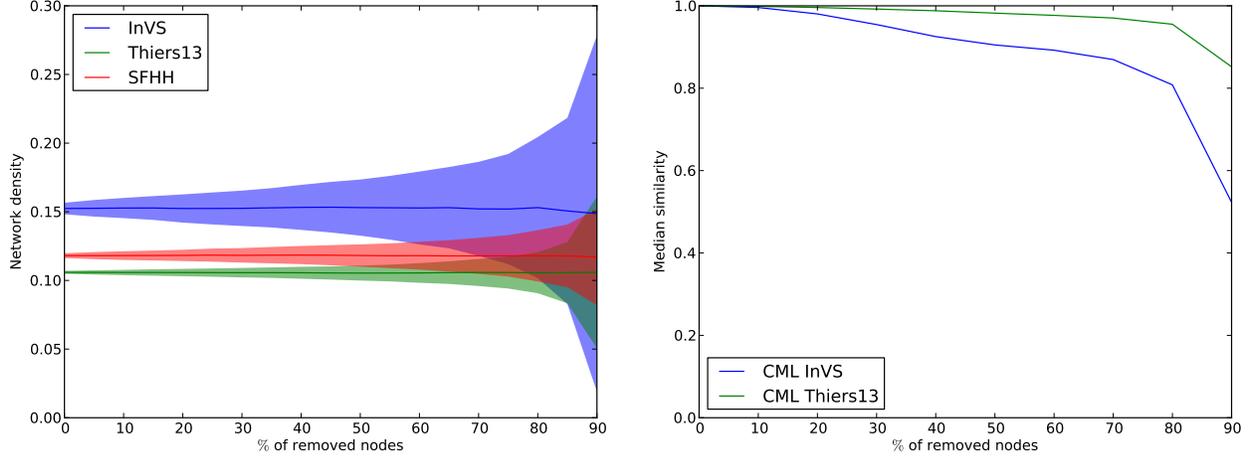


FIG. S2. **Effect of sampling on network density and on the similarity of contact matrices.** (Left) Density ρ of the aggregated network of contacts as a function of the fraction f of nodes excluded. The shaded areas represent mean $\rho \pm$ s.e.m.. (Right) Median cosine similarities between the link density contact matrices (CML) of resampled and full data sets, as a function of f , for the structured populations (high school and offices). Results are averaged, for each value of f , over 1,000 realisations for the density and over 100 realisations for the similarities.

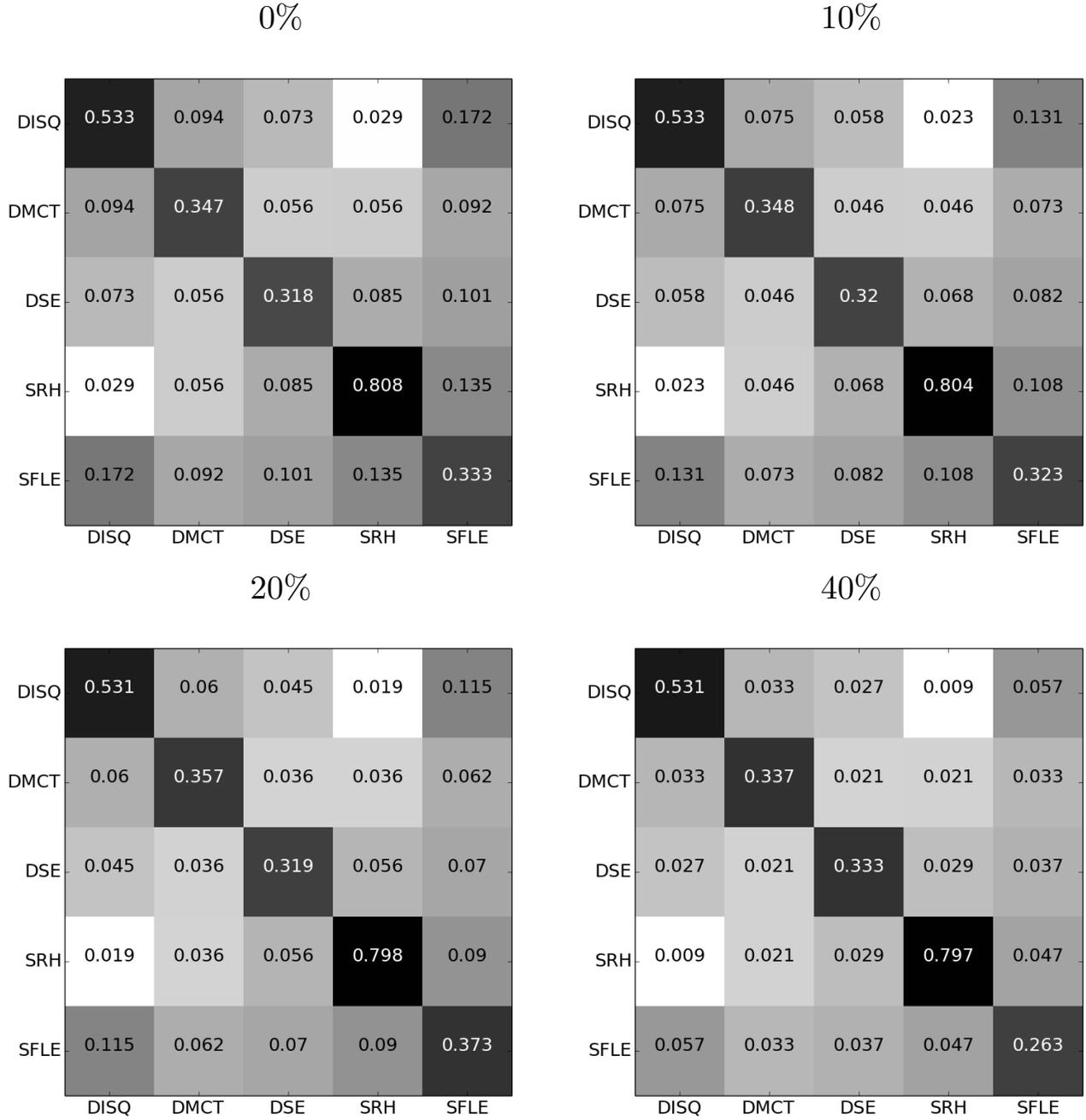


FIG. S3. **Effect of sampling: link density contact matrices (InVS).** Comparison of link density contact matrices for the office building, for different fractions of excluded nodes, f , with the original one ($f = 0$). Each matrix element AB gives the number of links between nodes of department A and nodes of department B in the contact network, normalised by the maximum possible number of such links. For each value of f , each matrix element is an average over 100 realisations of the sampling.

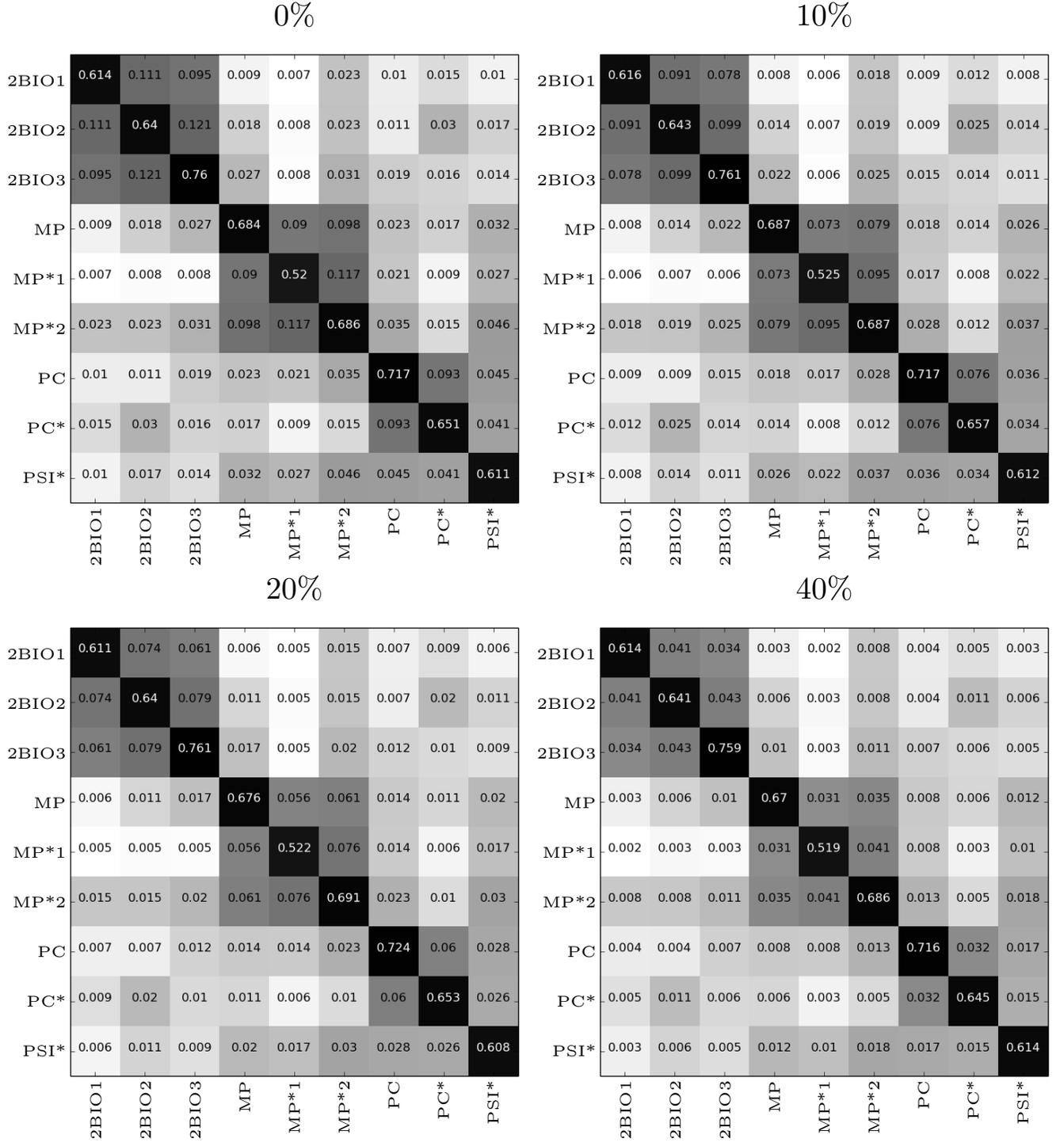


FIG. S4. **Effect of sampling: link density contact matrices (Thiers13).** Comparison of the link density contact matrices for the high school, for different fractions f of excluded nodes, with the original one ($f = 0$). Each matrix element AB gives the number of links between nodes of class A and nodes of class B in the contact network, normalised by the maximum possible number of such links. For each value of f , each matrix element is an average over 100 realisations of the sampling.

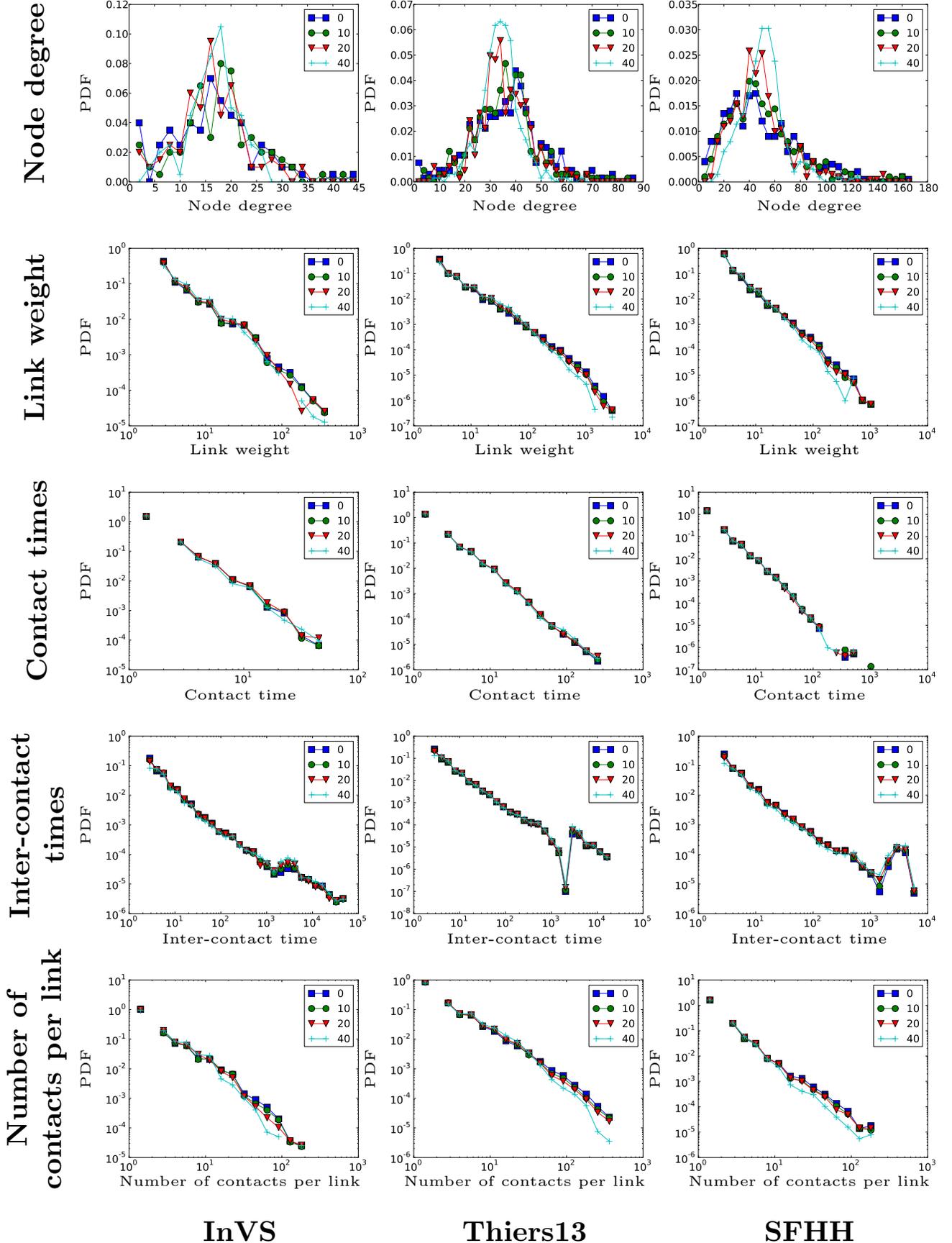


FIG. S5. **Properties of the reconstructed contact network.** Same as Fig. S1 but for the reconstructed networks: Distributions of structural (degrees and weights in the aggregated contact network) and temporal (contact times, inter-contact times, number of contacts per link) properties of the surrogate contact networks, for different fractions f of nodes excluded. For each value of f , the distributions are computed on a single reconstructed network.

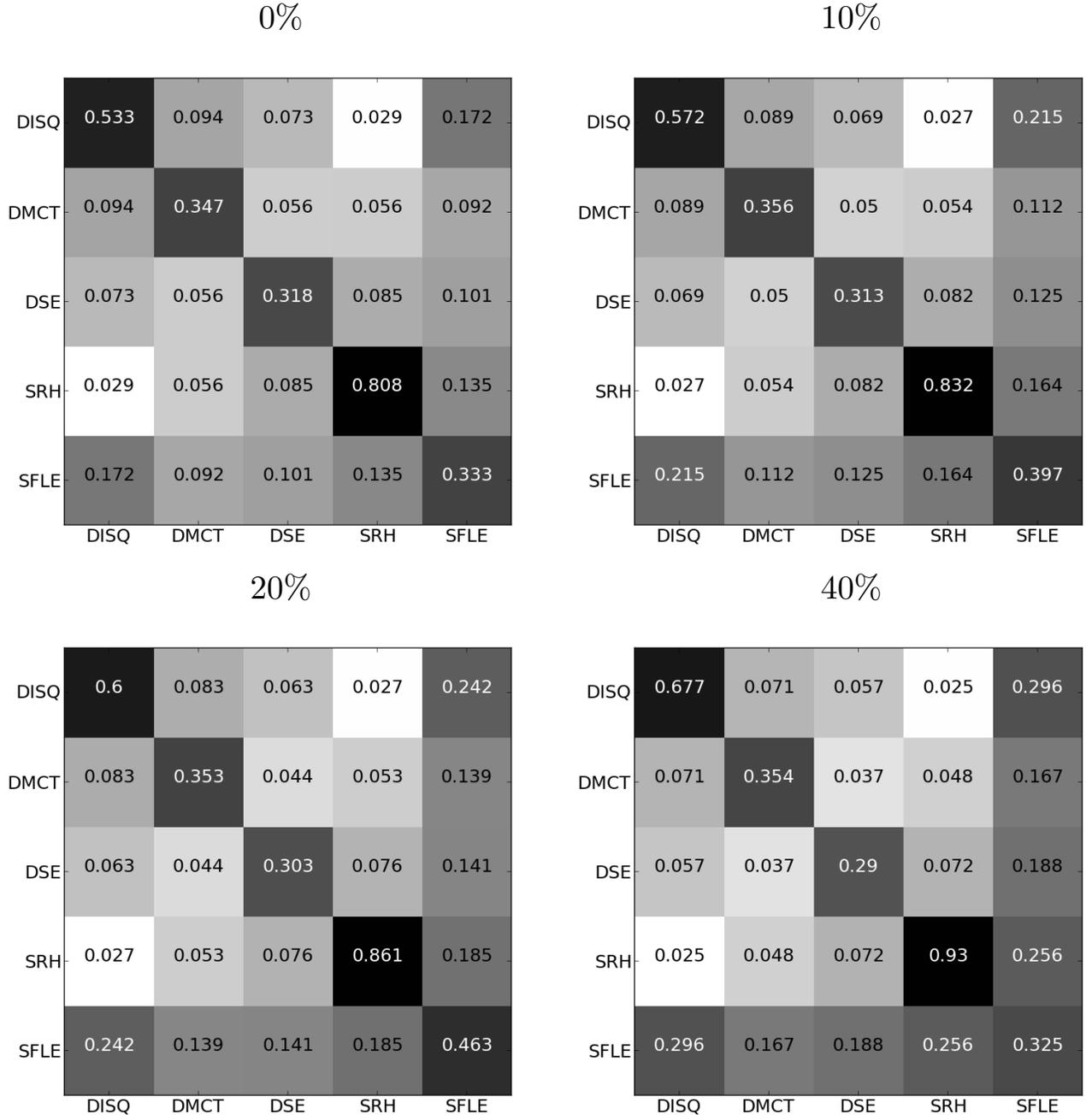


FIG. S6. **Properties of the reconstructed contact network: link density contact matrices (InVS).** Comparison of link density contact matrices for the reconstructed network of the office building data, for different values of the fraction f of excluded nodes, with the original one ($f = 0$). For each value of f , each matrix element is an average over 100 realisations of the sampling.

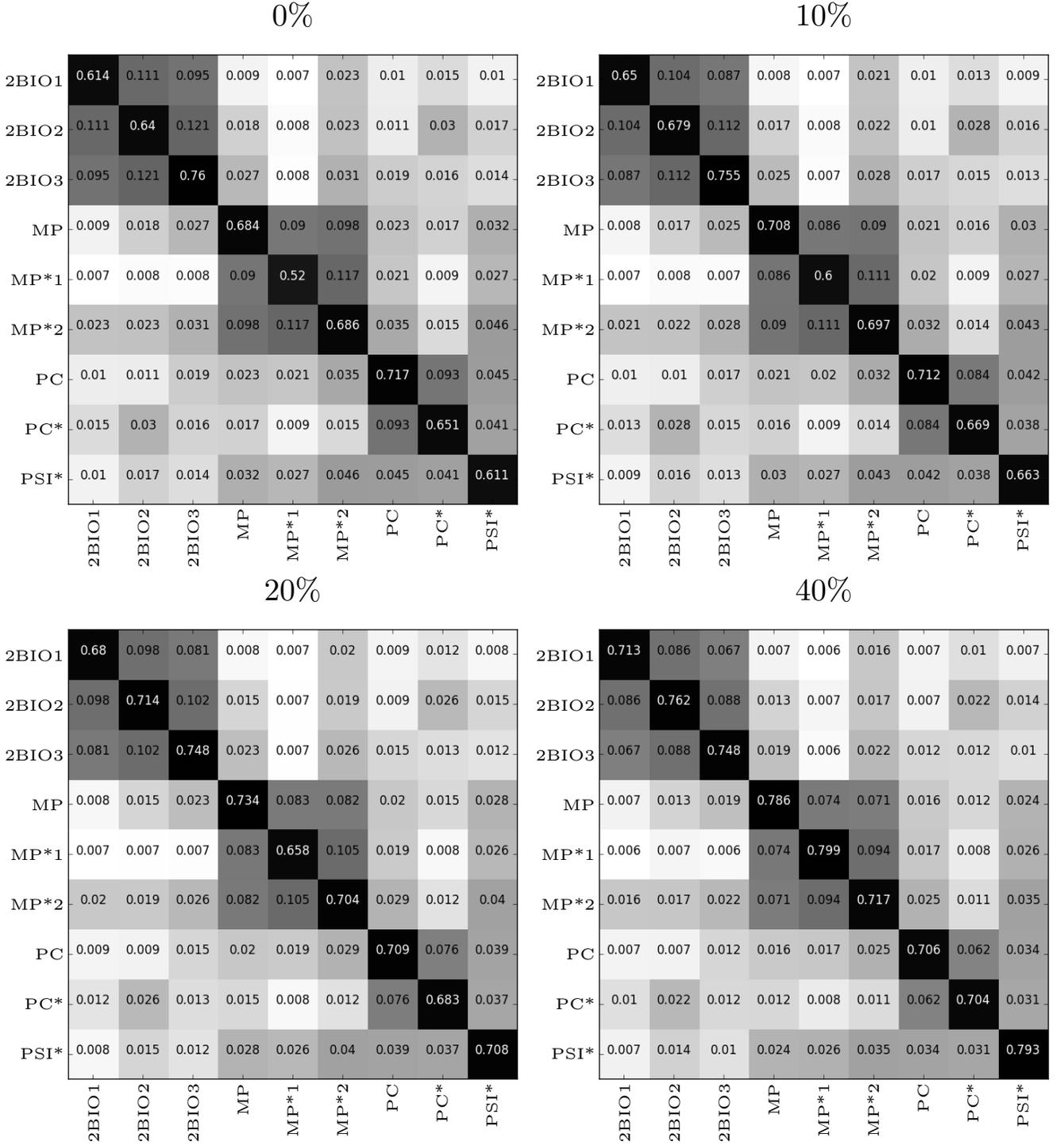


FIG. S7. **Properties of the reconstructed contact network: link density contact matrices (Thiers13).** Comparison of link density contact matrices for the reconstructed network of the high school data, for different values of the fraction f of excluded nodes, with the original one ($f = 0$). For each value of f , each matrix element is an average over 100 realisations of the sampling.

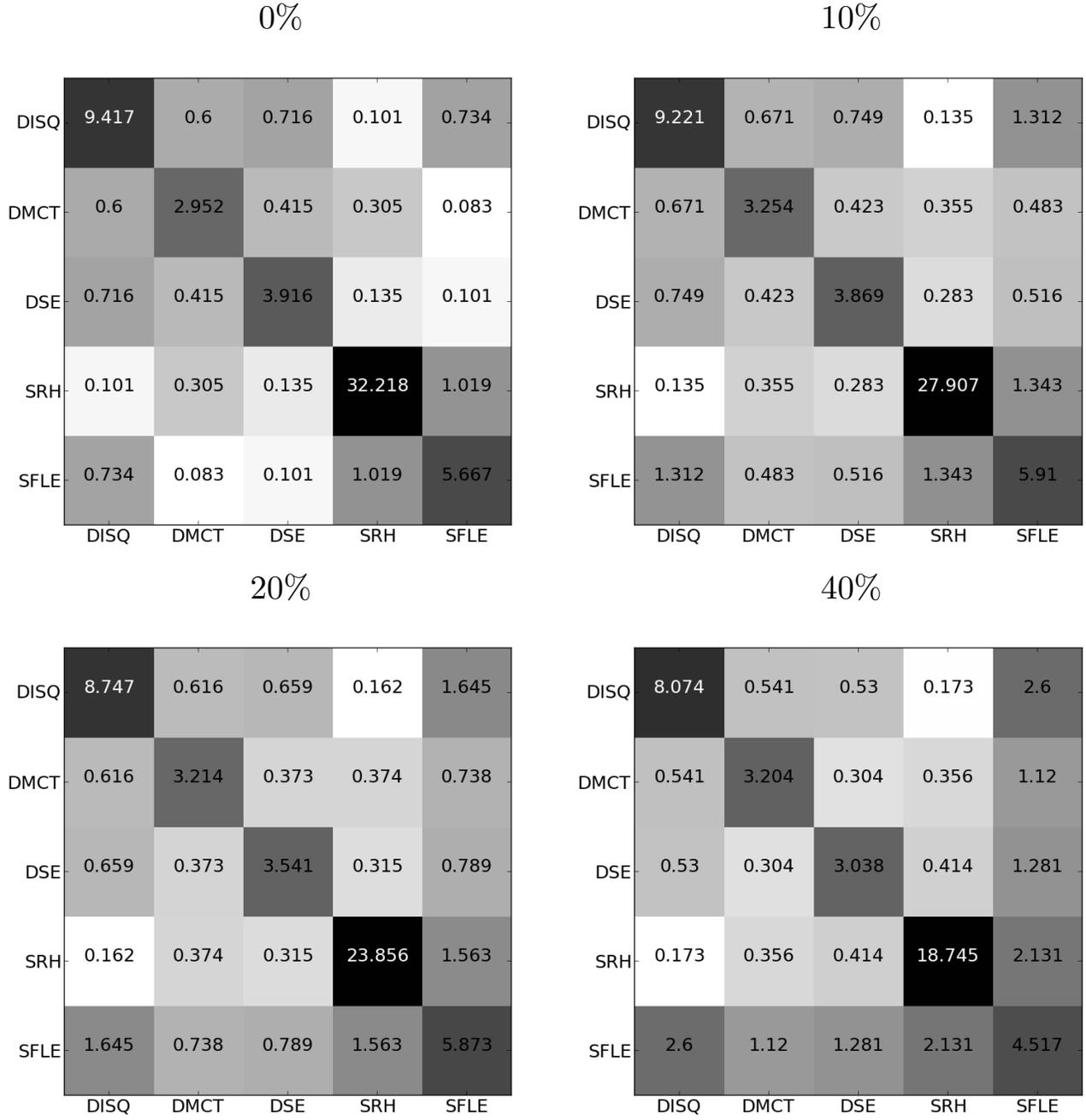


FIG. S8. **Properties of the reconstructed contact network: time density contact matrices (InVS).** Comparison of the contact time density contact matrices for the reconstructed network of the office building data, for different fractions of excluded nodes, f , with the original one ($f = 0$). Each matrix element AB gives the average time spent in contact between a node of department A and a node of department B . For each value of f , each matrix element is an average over 100 realisations of the sampling.

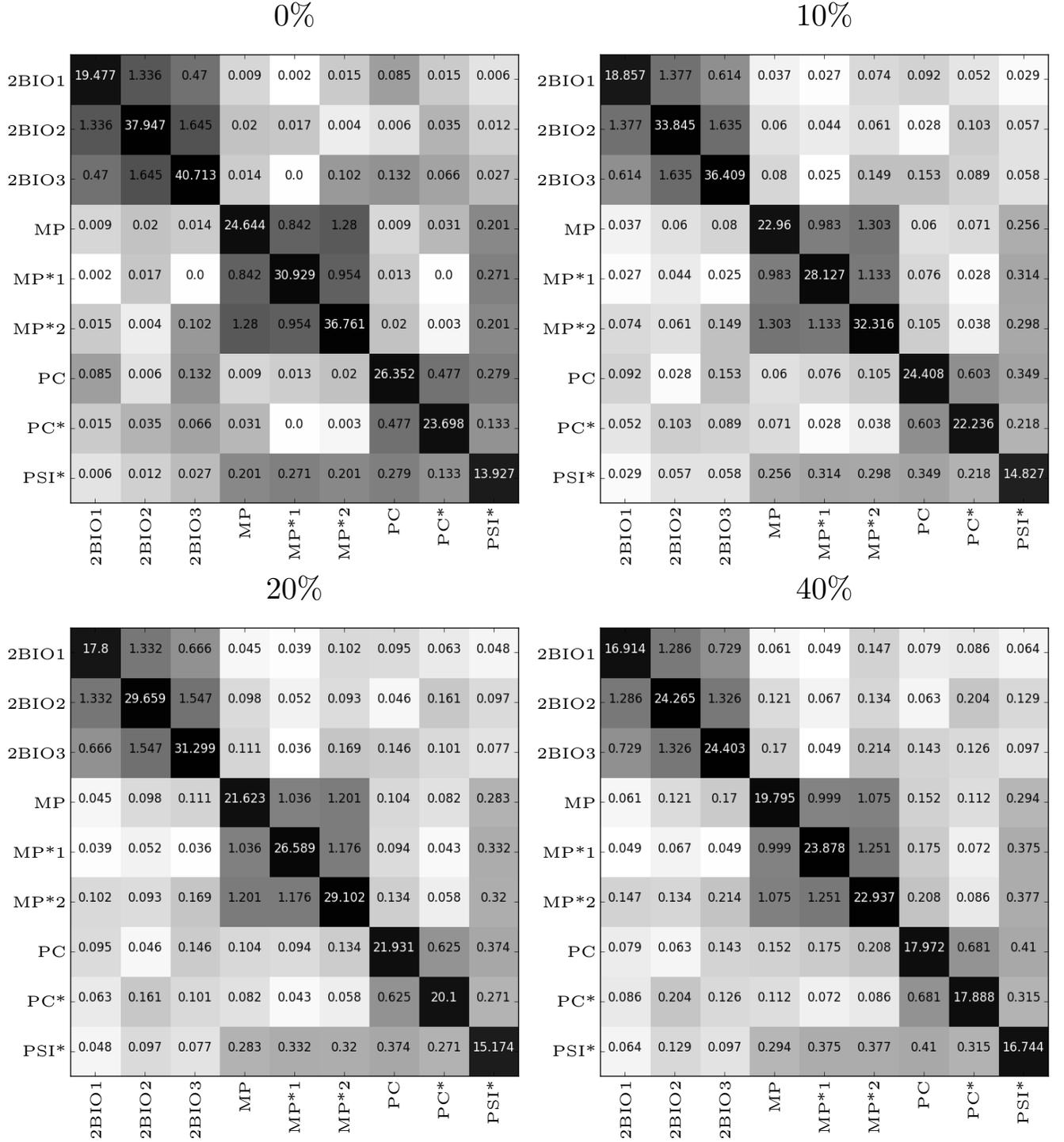


FIG. S9. **Properties of the reconstructed contact network: time density contact matrices (Thiers13).** Contact time density contact matrices for the reconstructed network of the high school data, for different fractions of nodes excluded, f . Each matrix element AB gives the average time spent in contact between a node of class A and a node of class B . For each value of f , each matrix element is an average over 100 realisations of the sampling.

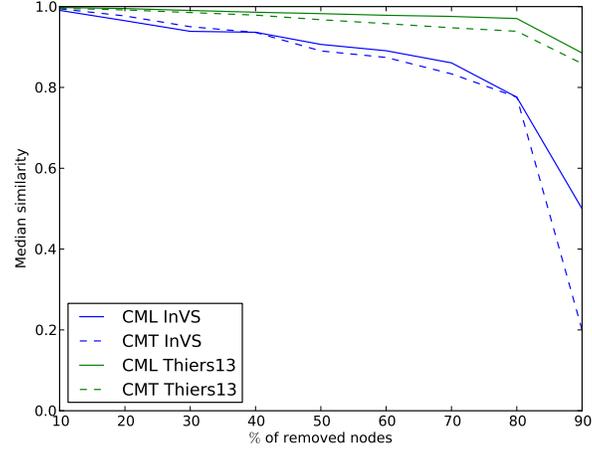


FIG. S10. **Similarity of contact matrices under sampling and reconstruction.** Median cosine similarity between the link density (CML) and contact time density (CMT) contact matrices computed for the reconstructed network and for the original contact matrices, as a function of the fraction f of removed nodes. For each value of f , the median is computed over 100 realisations of the reconstruction.

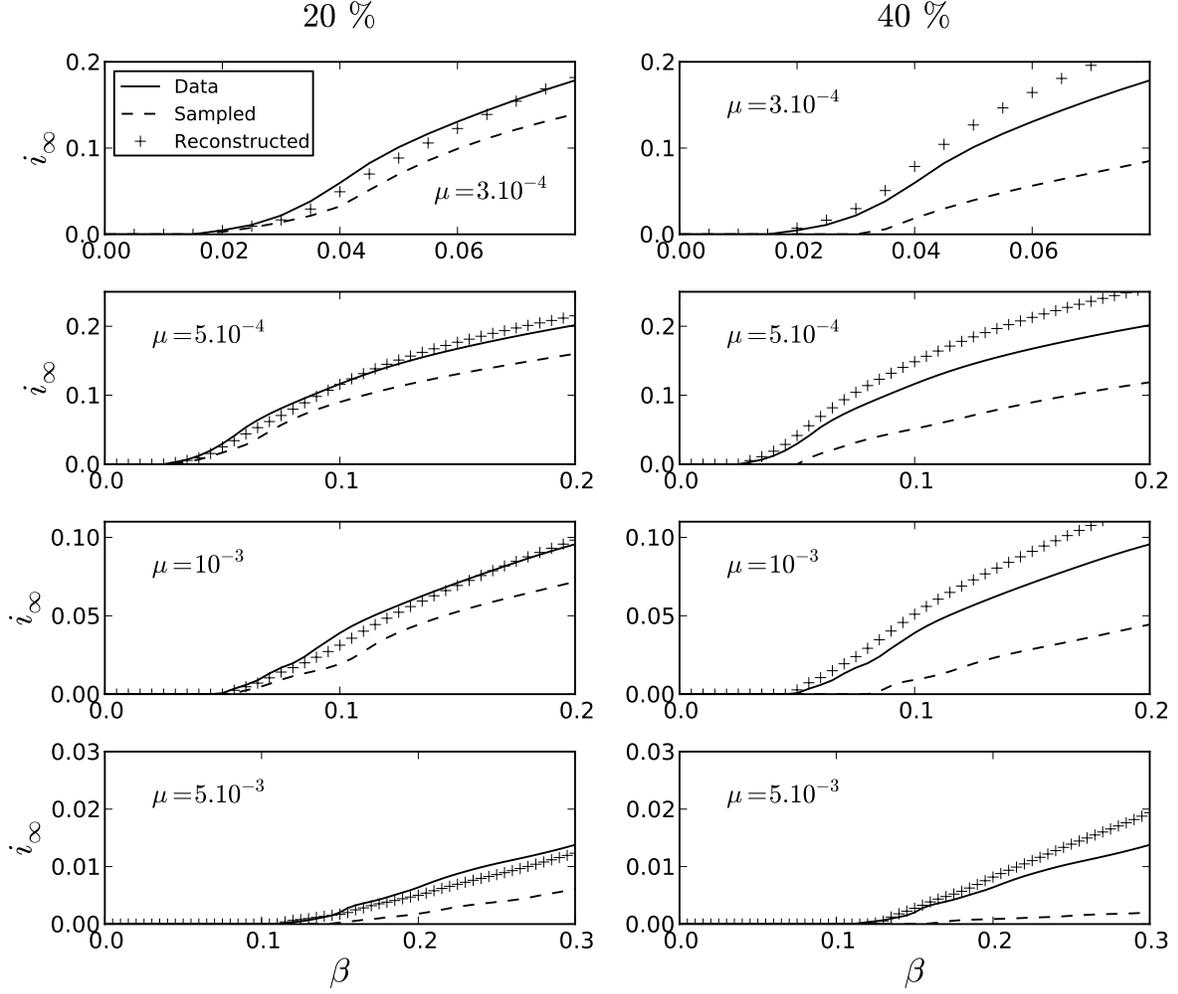


FIG. S11. **Phase diagram of the SIS model for original, resampled and reconstructed contact networks (InVS data set).** Each panel shows the stationary value i_∞ of the prevalence in the stationary state of the SIS model, computed as described in the Methods section, as a function of β , for several values of μ . Here we consider the example of the InVS data set. The epidemic threshold corresponds to the transition between $i_\infty = 0$ and $i_\infty > 0$. The prevalence curves are computed in each case using either the whole data set (continuous lines), resampled data (dashed lines) or reconstructed contact networks (pluses). The fraction of excluded nodes in the resampling is $f = 20\%$ for the left column and $f = 40\%$ for the right column.

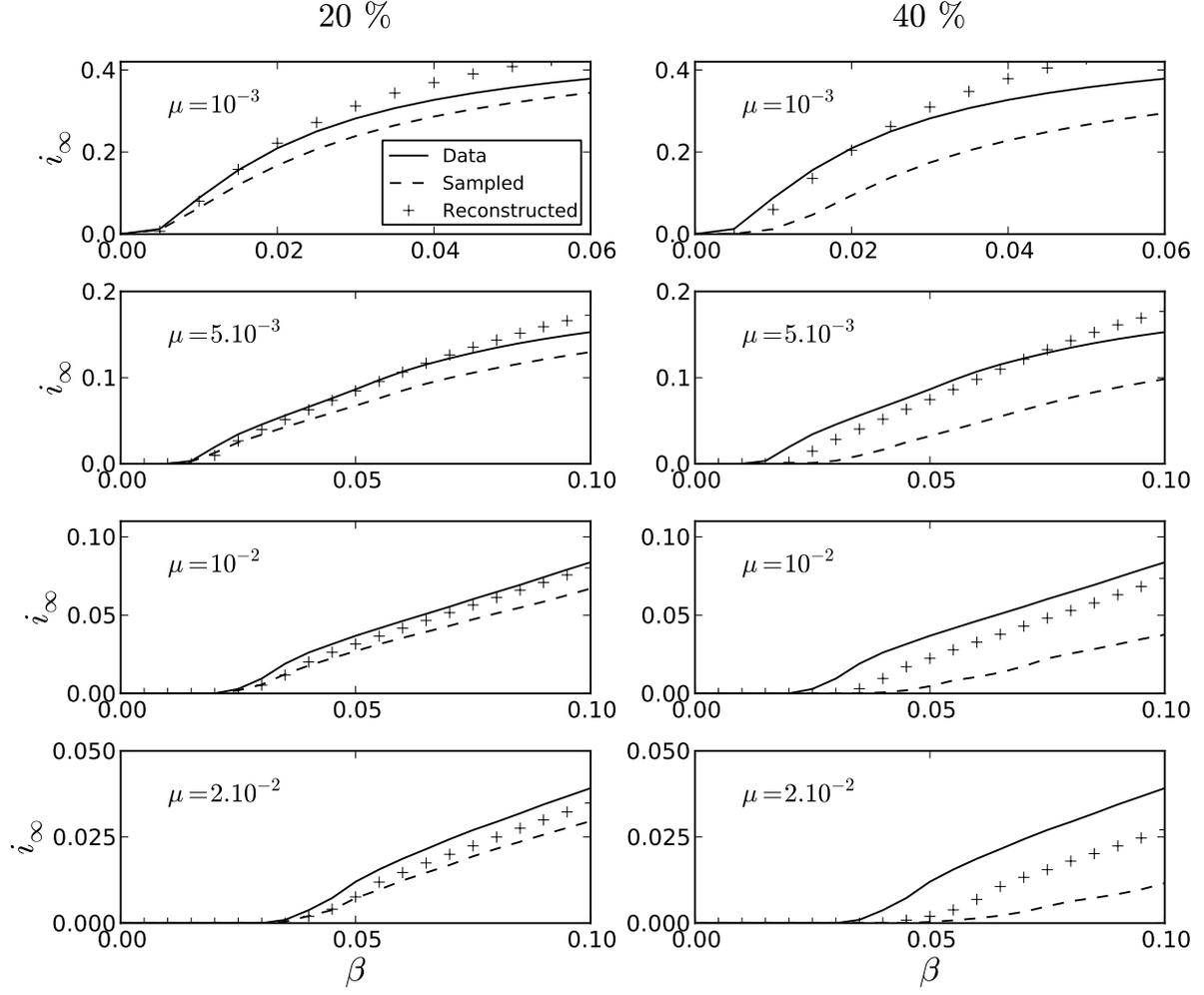


FIG. S12. Phase diagram of the SIS model for original, resampled and reconstructed contact networks (SFHH data set). Same as Fig. S11 for the SFHH (conference) data set.

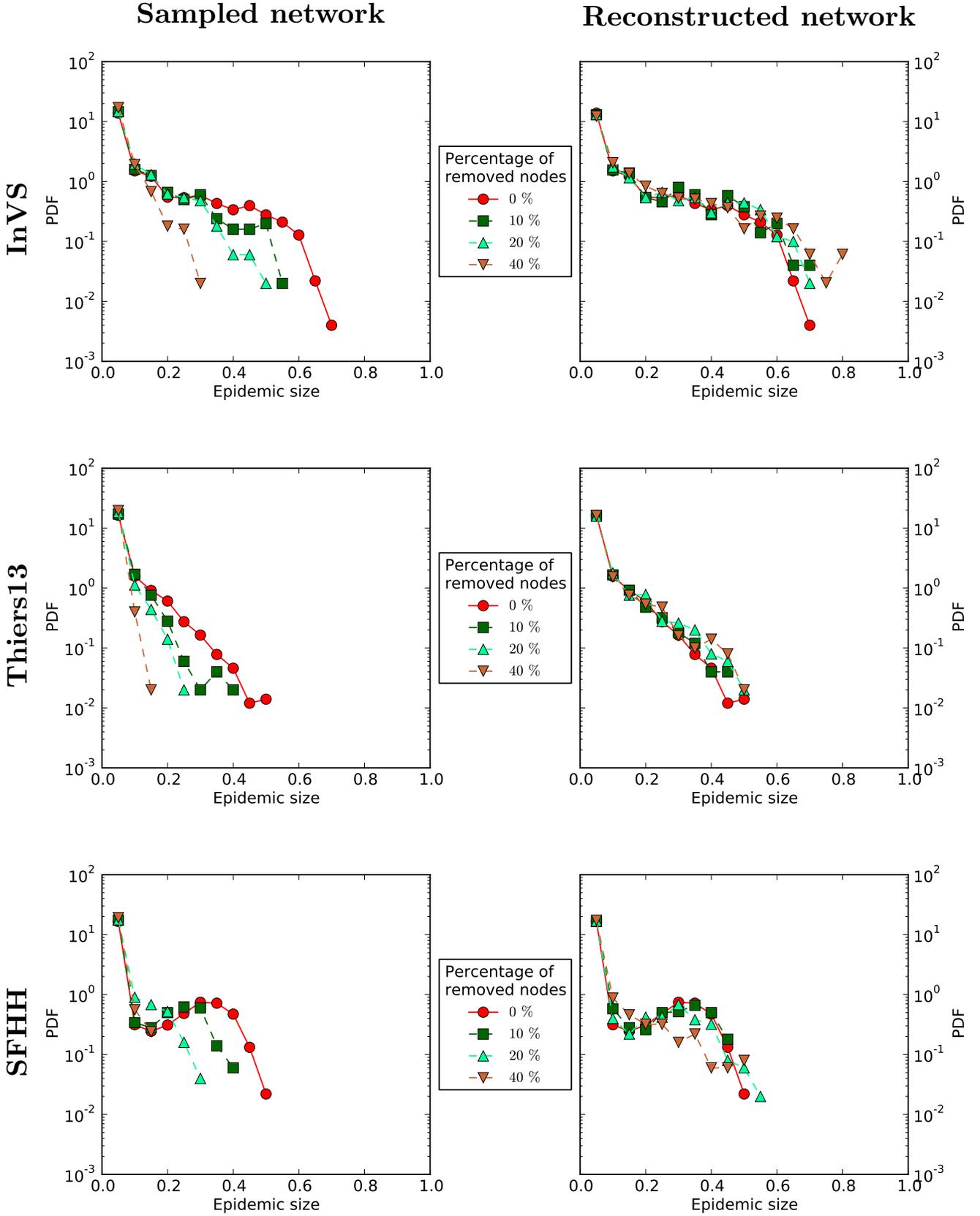


FIG. S13. **Outcome of SIR epidemic simulations on reconstructed networks for different parameter values.** Distribution of epidemic sizes for simulations of SIR processes on reconstructed networks and on the whole data set (case $f = 0$), for different values of the fraction f of excluded nodes. Each distribution is computed on 1,000 simulations of the SIR process. Here $\beta = 0.004$, $\beta/\mu = 500$ for the InVS, $\beta/\mu = 50$ for Thiers13 and $\beta/\mu = 30$ for SFHH. Each distribution is computed on 1,000 simulations of the SIR process.

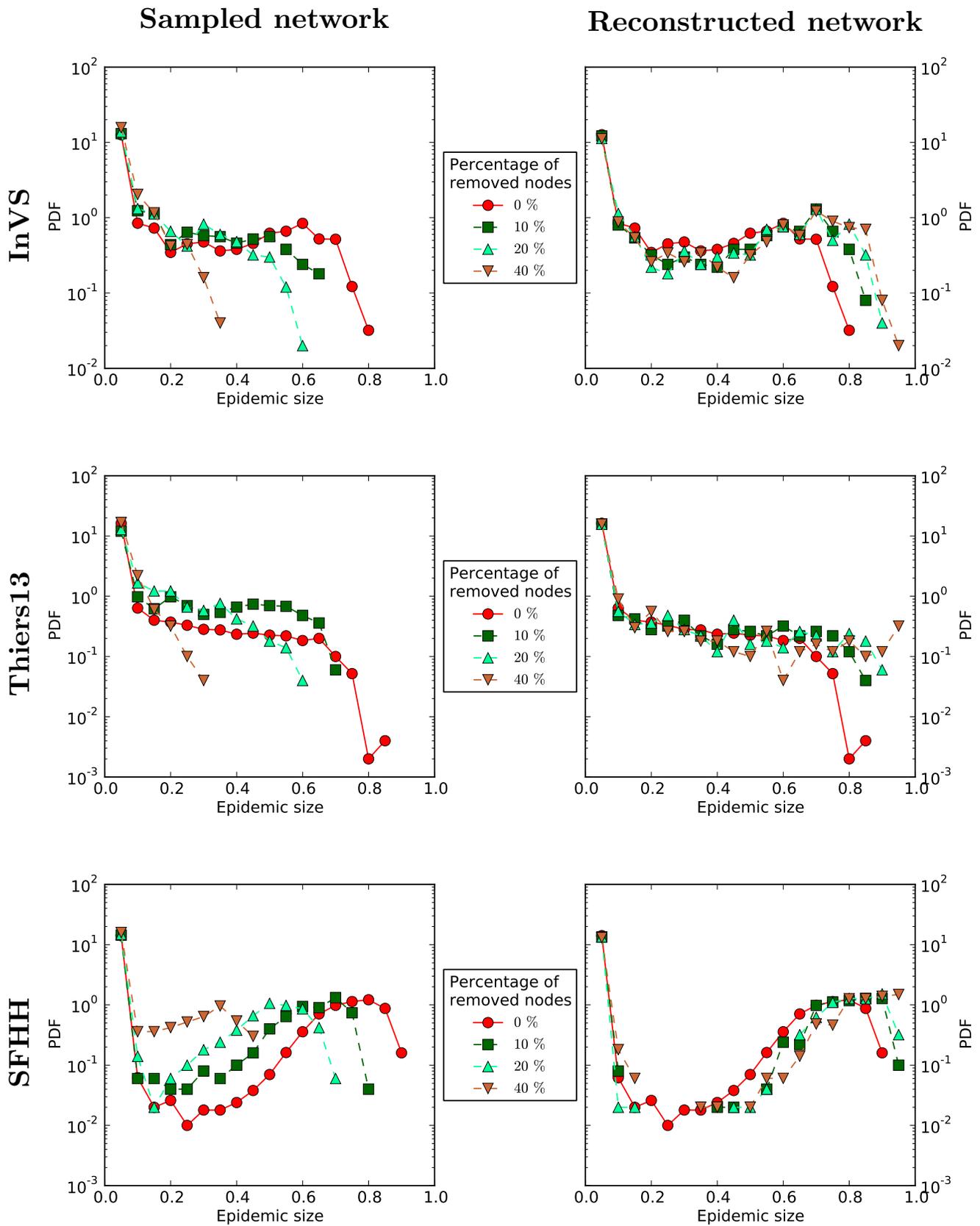


FIG. S14. Outcome of SIR epidemic simulations on reconstructed networks for different parameter values. Same as Fig. S13 for $\beta = 0.004$, $\beta/\mu = 1000$ (InVS) or $\beta/\mu = 100$ (Thiers13 and SFHH). Each distribution is computed on 1,000 simulations of the SIR process.

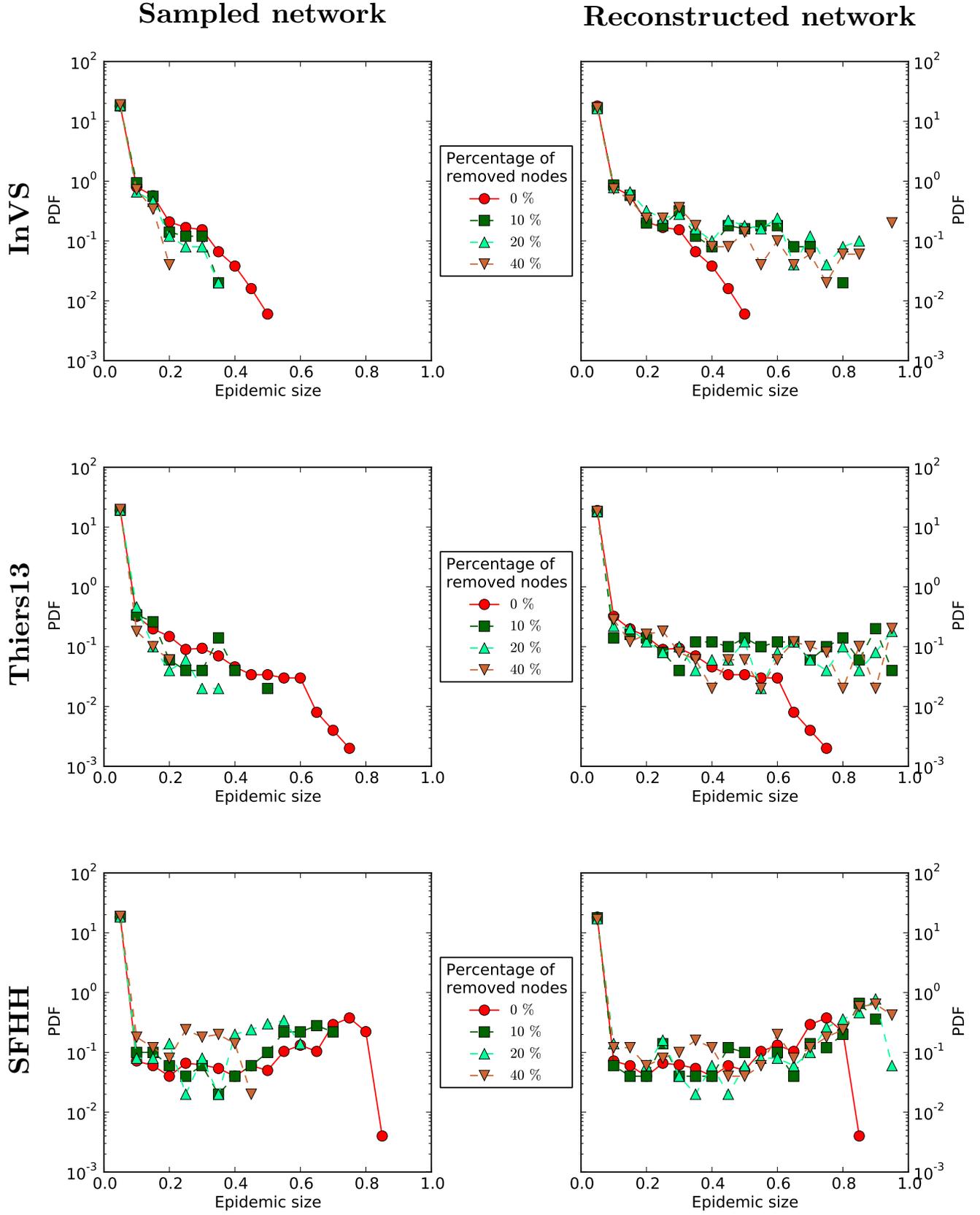


FIG. S15. Outcome of SIR epidemic simulations on reconstructed networks for different parameter values. Same as Fig. S13 for $\beta = 0.04$, $\beta/\mu = 1000$ (InVS) or $\beta/\mu = 100$ (Thiers13 and SFHH). Each distribution is computed on 1,000 simulations of the SIR process.

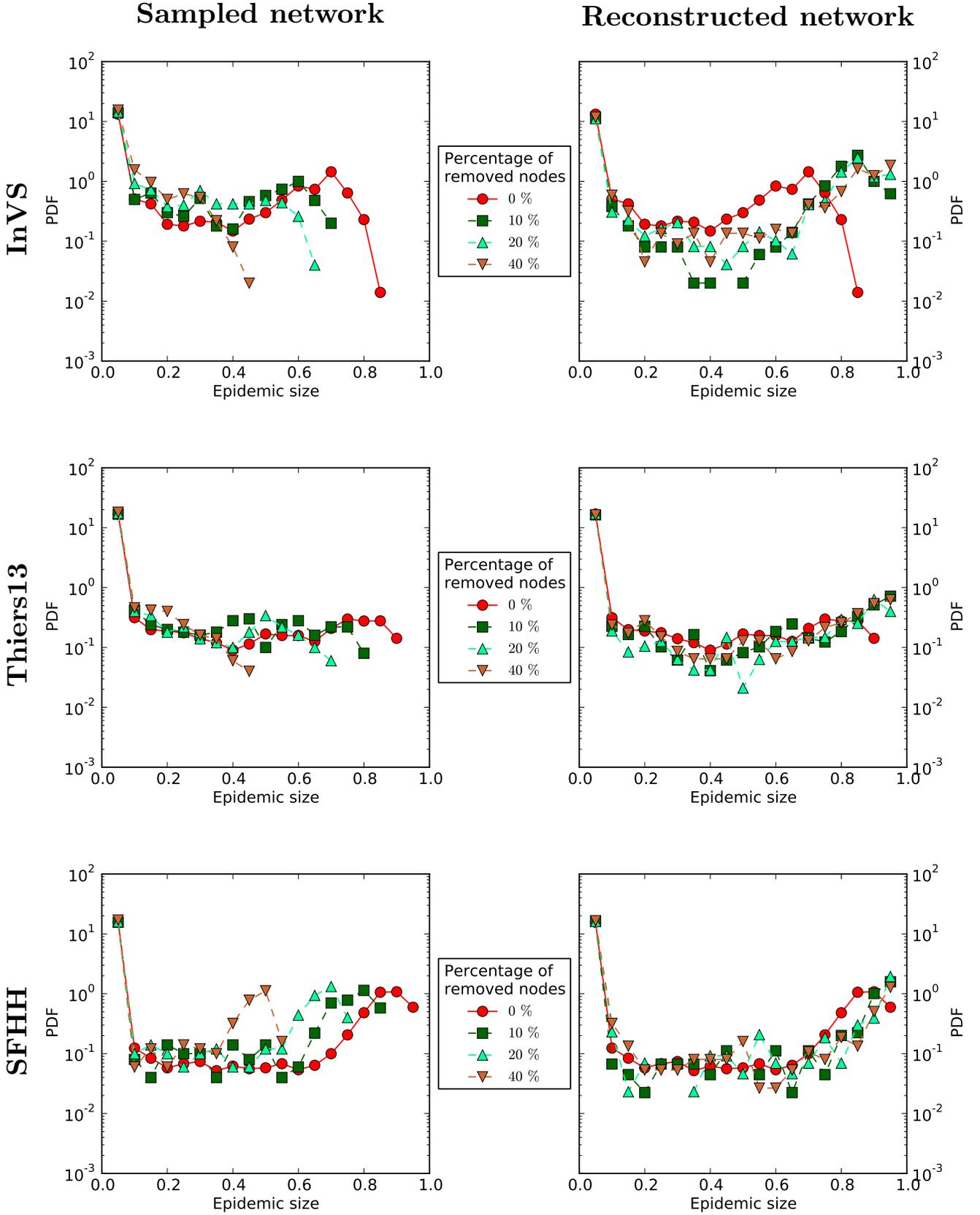


FIG. S16. Outcome of SIR epidemic simulations on reconstructed networks for different parameter values. Same as Fig. S13 for $\beta = 0.04$, $\beta/\mu = 4000$ (InVS) or $\beta/\mu = 400$ (Thiers13 and SFHH).