

A fast and accurate algorithm for inferring sparse Ising models via parameters activation to maximize the pseudo-likelihood

Silvio Franz,^{1,2} Federico Ricci-Tersenghi,^{2,3} and Jacopo Rocchi¹

¹*LPTMS, Université Paris-Sud 11, UMR 8626 CNRS, Bât. 100, 91405 Orsay Cedex, France*

²*Dipartimento di Fisica Università, La Sapienza, Piazzale Aldo Moro 5, I-00185 Roma, Italy*

³*INFN-Sezione di Roma 1, and CNR-Nanotec, Rome unit, P.le A. Moro 5, I-00185, Roma, Italy*
(Dated: February 19, 2019)

We propose a new algorithm to learn the network of the interactions of pairwise Ising models. The algorithm is based on the pseudo-likelihood method (PLM), that has already been proven to efficiently solve the problem in a large variety of cases. Our present implementation is particularly suitable to address the case of sparse underlying topologies and it is based on a careful search of the most important parameters in their high dimensional space. We call this algorithm Parameters Activation to Maximize Pseudo-Likelihood (PAMPL). Numerical tests have been performed on a wide class of models such as random graphs and finite dimensional lattices with different type of couplings, both ferromagnetic and spin glasses. These tests show that PAMPL improves the performances of the fastest existing algorithms.

The Ising model is a graphical model whose parameters $\{J_{ij}, h_i\}$ can be tuned in order to describe stationary distributions of binary variables, s_i , according to the weight $P(\underline{s}) \sim \exp\left(\sum_{i<j} J_{ij}s_i s_j + \sum_i h_i s_i\right)$. In many practical problems in different domains - e.g. physics, biology, neuroscience, finance, sociology - the topology of the graph and the values of the couplings are unknown and they need to be reconstructed from the data. The inverse Ising problem aims to find the parameters of the model that best fit the data.

From the original attempt to solve this problem [1], many techniques of statistical mechanics and machine learning have been developed [2–12] to study different cases. The need to develop approximate algorithms can be understood from the observation that the likelihood depends on the partition function, which is generally intractable. Among these methods, the pseudo-likelihood [13] has been proven to be particularly efficient, leading to polynomial algorithms which give the exact solution in the limit of infinite sampling. Methods based on the pseudo-likelihood need to be complemented with a threshold procedure, implemented a posteriori or through a regularization function. An improvement of this method based on a decimation scheme was presented in [11]. The decimation based algorithm has been shown to outperform existing algorithms based on the pseudo-likelihood method in terms of the quality of the reconstructed graph and it has been commonly used in a variety of contexts [14–16]. Our aim is to improve this algorithm in the case of sparse graphs. In fact, in this case, the underlying structure is closer to an empty graph than to a fully connected graph and we would like to avoid to explore the full set of parameters in the inference process, while maintaining the same quality in the inferred graph. While the decimation step is $O(N^2)$, our elementary operation is $O(N)$. We begin formulating the Inverse Ising problem. We discuss the pseudo-likelihood

method and the present implementation. Finally we describe the results of our algorithm in a wide class of Ising models with a comparison with the fast Minimum Probability Flow (MPF) [4], showing that the two methods have similar execution times and that ours outperform the other in terms of the quality of the reconstructed graphs.

An Ising model in the absence of local fields is defined by the Hamiltonian $H(\underline{s}) = -\sum_{i<j} J_{ij}s_i s_j$. After the observation of M independent equilibrium configurations, the problem of inferring the couplings can be formulated in terms of the Bayes theorem $P(J|\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}) \propto P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}|J)P(J)$, where the two functions on the r.h.s. are named likelihood and prior, respectively. If we assume to be in a Bayes optimal case setting where we do not need to introduce local fields in the model, the log-likelihood function is defined by $\mathcal{L}(J) = M^{-1} \log P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}|J)$ and reads

$$\mathcal{L}(J) = \frac{\beta}{M} \sum_{\mu=1}^M \sum_{i<j} J_{ij} s_i^{(\mu)} s_j^{(\mu)} - \log Z(J), \quad (1)$$

where $Z(J) = \sum_{\underline{s}} e^{\beta \sum_{i<j} J_{ij} s_i s_j}$ is the partition function of the problem, and β an inverse temperature. Optimizing $\mathcal{L}(J)$ over the parameters of the model leads to

$$\frac{\partial \mathcal{L}(J)}{\partial J_{ij}} = \beta \left(\langle s_i s_j \rangle_{\text{Data}} - \langle s_i s_j \rangle_{\text{Model}} \right), \quad (2)$$

where we defined $M \langle s_i s_j \rangle_{\text{Data}} = \sum_{\mu=1}^M s_i^{(\mu)} s_j^{(\mu)}$ and $\beta \langle s_i s_j \rangle_{\text{Model}} = \partial_{J_{ij}} \log Z$. The J that maximizes the log-likelihood is such that

$$\langle s_i s_j \rangle_{\text{Data}} = \langle s_i s_j \rangle_{\text{Model}}. \quad (3)$$

This formulation involves the computation of the partition function, which is a complicated object. The log-pseudo-likelihood [17] is introduced to deal with this dif-

ficulty. It is defined by

$$\mathcal{S}(J) = \frac{1}{M} \sum_{r=1}^N \sum_{\mu=1}^M \log p(s_r^{(\mu)} | \underline{s}_{\setminus r}^{(\mu)}) \quad (4)$$

where $p(s_r | \underline{s}_{\setminus r}) = \left[1 + e^{-2\beta s_r \sum_{j \neq r} J_{rj} s_j} \right]^{-1}$, and as discussed in the Appendix A, maximizing \mathcal{S} leads to the correct solution in the infinite sampling limit.

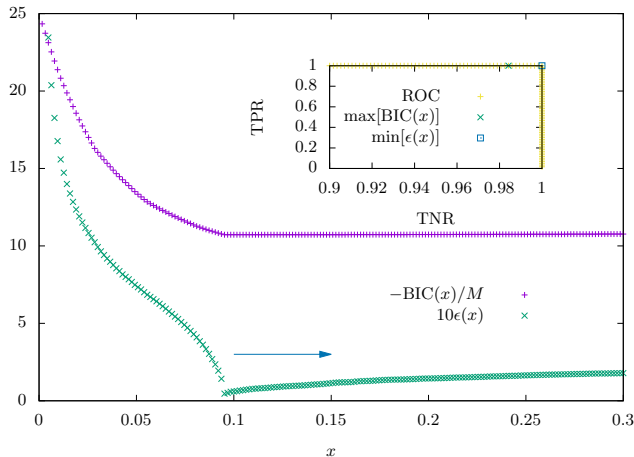


FIG. 1. 2-D ferromagnetic lattice with $N = 36$ and free boundary conditions, $M = 5000$, $\beta = 0.5$. Main Figure: (a) $-\text{BIC}/M$ and 10ϵ as a function of x , fraction of couplings, for the activation-based method. x is 0 at the beginning and it increases as the iteration proceeds. Only the first part of the iteration is presented. The Arrow specifies the directions in which the iteration moves. Inset: ROC curve, with the TNR on the x-axis and the TPR on the y-axis.

The use of the pseudo-likelihood has been already shown to be very useful in the context of the Inverse Ising problem [6, 13]. The standard implementation of this method consists in maximizing each of the N local likelihood functions separately, thus getting two different estimates for each coupling. When complemented with a post-optimization parameter thresholding procedure, this method can be shown to reconstruct arbitrary Ising models [12]. Anyway, this scheme relies on the delicate choice of the threshold and leads to estimated couplings that are systematically smaller than the true values. This problem was first addressed in [11] to eliminate the bias in the coupling estimation, where an iterative decimation based approach was developed. A maximization over the pseudo-likelihood is alternated with a decimation step where the smallest estimated couplings are set to zero. More details and examples are provided in Appendix B.

Here we propose an improvement of this algorithm especially suitable for sparse graphs. In fact, in this case, starting from the complete graph and decimating couplings requires a long time before reaching the correct stopping point. On the contrary, it would be wiser to

have an iterative algorithm that starts from the empty graph, and add links sequentially. We call this algorithm Parameters Activation to Maximize Pseudo-Likelihood (PAMPL). In order to add the correct links, we search for the directions, in the parameter space, that give the largest gain in the log-pseudo-likelihood. The change in \mathcal{S} due to a change in the coupling J_{ij} is estimated using a second order approximation,

$$\mathcal{S}(J + \Delta J_{ij}) = \mathcal{S}(J) + \mathcal{S}'(J) \Delta J_{ij} + \frac{\mathcal{S}''(J)}{2} \Delta J_{ij}^2 + \dots, \quad (5)$$

where prime denotes differentiation with respect to J_{ij} . Couplings are updated with one step of the Newton method,

$$0 = \mathcal{S}'(J_{ij}) + \mathcal{S}''(J_{ij}) \Delta J_{ij}, \quad (6)$$

and ranked in an ascending order according to the values of the quantities

$$\Delta \mathcal{S}_{ij} = -\frac{1}{2} \frac{\mathcal{S}'^2(J_{ij})}{\mathcal{S}''(J_{ij})}, \quad (7)$$

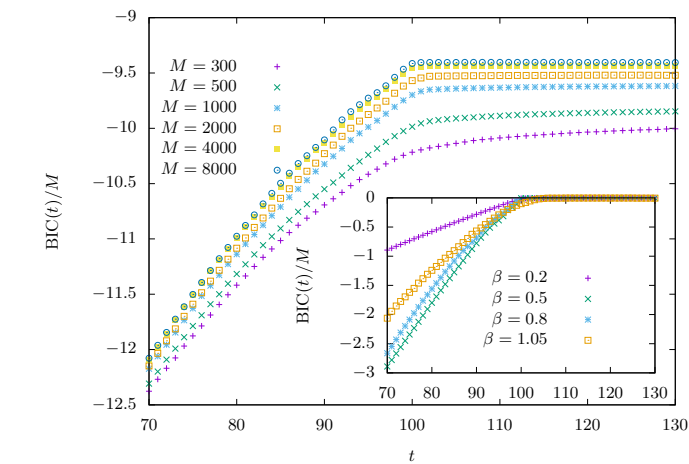
obtained by plugging eq. (6) in eq. (5). Finally, the first K are included in the set of non-zero couplings \mathbb{J} and optimized as explained below. This procedure is iterated adding more and more couplings in \mathbb{J} at each iteration. In order to keep this elementary step $O(N)$, updating and sorting need to be done carefully. In particular, there is no need to update all of the $\Delta \mathcal{S}_{ij}$ at each step, since only $O(N)$ of them are affected by the presence of a new coupling in \mathbb{J} . Moreover, since most of them are small, we don't need to order $O(N^2)$ elements, but only $O(N)$. We use the Bayesian Information Criterion (BIC), introduced in [18], to locate the stopping point of the iteration. For our purposes, it is defined by

$$\text{BIC} = 2MS^* - k \log M, \quad (8)$$

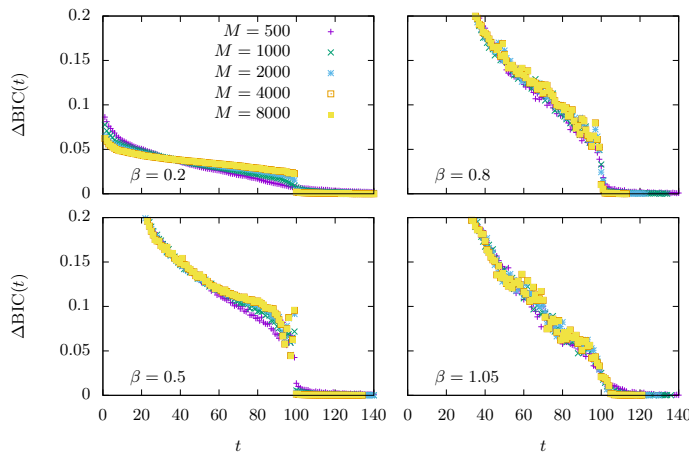
where $k = ||\mathbb{J}||$ is the cardinality of the set \mathbb{J} , corresponding to the number of links used to describe the observations, and \mathcal{S}^* is the maximum of \mathcal{S} found optimizing the couplings in \mathbb{J} . In order to perform the optimization step on the couplings of \mathbb{J} , we used the LBFGS [19] algorithm and a simple gradient ascent. Results obtained in the two cases are the same within numerical errors and, since the second one is faster, it is particularly appropriate for large systems. In the following we update one coupling per iteration time. More details about the algorithm are provided in Appendix C.

We generate independent equilibrium configurations from given graphs using a Monte Carlo sampling algorithm. Then, during the inference process, we compare the inferred graph with the original one using the measure

$$\epsilon = \sqrt{\frac{\sum_{i < j} (J_{ij} - J_{ij}^*)^2}{\sum_{i < j} J_{ij}^2}}, \quad (9)$$



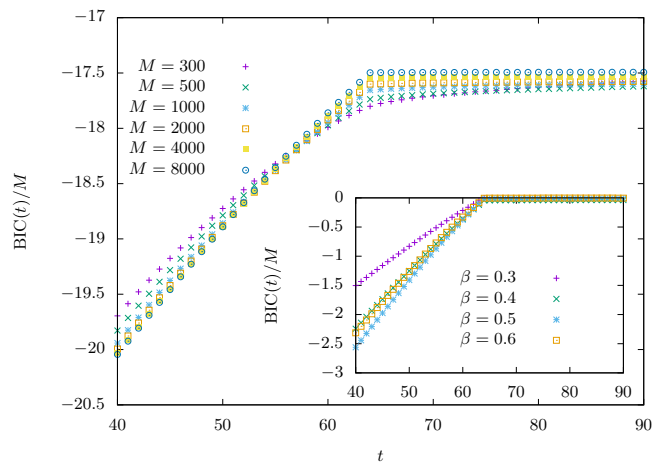
(a)



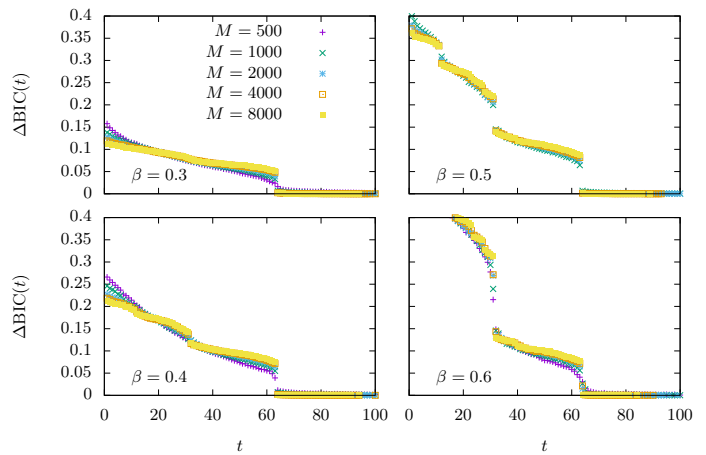
(b)

FIG. 2. Random regular spin glass with $N = 50$ spins, $c = 4$, 100 couplings, $\beta_c = 0.524$. (a): Main Figure: average of BIC/M as a function of the iteration time over 10 runs of the activation algorithm, for different values of M at $\beta = 0.8$. Inset: average of BIC/M as a function the iteration time over 10 runs of the activation algorithm, for different values of the inverse temperature β , at $M = 8000$. BIC corresponding to different β has been shifted in order to fit in the same figure. (b): Average of $\Delta BIC(t)$ over 10 runs of the activation algorithm, for different values of the inverse temperature β .

were J^* is the original set of couplings. We study graphs with ferromagnetic and spin glass interactions. We denote with spin glass graph systems with couplings equal to ± 1 with probability 0.5. We stress that no extra time is required to infer the structure of a spin glass topology compared to the ferromagnetic case. In Fig. 1 we show the behavior of our algorithm on a 2-D lattice with $N = 36$ and free boundary conditions. Inference is made after observing $M = 5000$ samples extracted at equilibrium at $\beta = 0.5$. In the inset we plot the ROC curves that give information on the fraction of true couplings



(a)



(b)

FIG. 3. 2D ferromagnetic diamond lattice with $N = 44$ spins, 64 couplings, $\beta_c = 0.609$. (a): Main Figure: average of BIC/M as a function of the iteration time over 10 runs of the activation algorithm, for different values of M at $\beta = 0.5$. Inset: average of BIC/M as a function the iteration time over 10 runs of the activation algorithm, for different values of the inverse temperature β at $M = 8000$. BIC corresponding to different β has been shifted in order to fit in the same figure. (b): Average of $\Delta BIC(t)$ over 10 runs of the activation algorithm, for different values of the inverse temperature β .

retrieved (true positive rate, TPR), and the fraction of non-existent couplings not created by the algorithm (true negative rate, TNR). Each point of the curve corresponds to the graph inferred at a particular stage of the iterative process. Ideally, for a perfect reconstruction, the inferred graph corresponds to the point $(1, 1)$. We observe that the maximum of the BIC does not coincide with the point where ϵ is minimum. This is due to the fact that after the activation of all the correct couplings of the graph, the BIC keeps growing for some other time steps before the penalty terms start to be effective. Despite this is-

sue, we notice that the correct stopping point is clearly recognizable by a visual inspection: this problem can be overcome easily, as will be shown below.

In Fig. 2-3 we study the performances of the algorithm with M and β for different topologies and sizes. As expected, inference becomes hard in the low temperature phase and when the dataset is too small. In fact, for small M the singular behavior of the BIC becomes smoother, and the detection of the stopping point is impossible. On the other hand, at larger values of β , more and more samples are required for a correct inference because most of the samples are very close to the ground state(s) and we lose information from fluctuations. We consider a Random Regular (RR) spin glass graph with $N = 50$ and $c = 4$ and a ferromagnetic diamond lattice of $N = 44$ spins. Diamond lattices [20], [21] are graphs constructed recursively from a single link corresponding to the generation $n = 0$. The generation $n = 1$ consists of 2 branches in parallel, each one made by 2 links in series. The generation $n = 2$ is obtained by applying the same transformation to the each link. The present case corresponds to the case $n = 3$. In this graphs there is a clear hierarchy between couplings and we show that the learning algorithm is clearly sensitive to it. In both cases we study the quantity $\Delta\text{BIC}(t) = [\text{BIC}(t) - \text{BIC}(t - 1)]/M$, averaged over 10 inference iterations, as a function of the iteration time. We observe that it becomes very small as soon as the correct graph structure is recovered, and that this threshold behavior is more evident in the vicinity of the phase transition. We find that a good stopping point can be defined when $\Delta\text{BIC} < 0.01$. The quality of

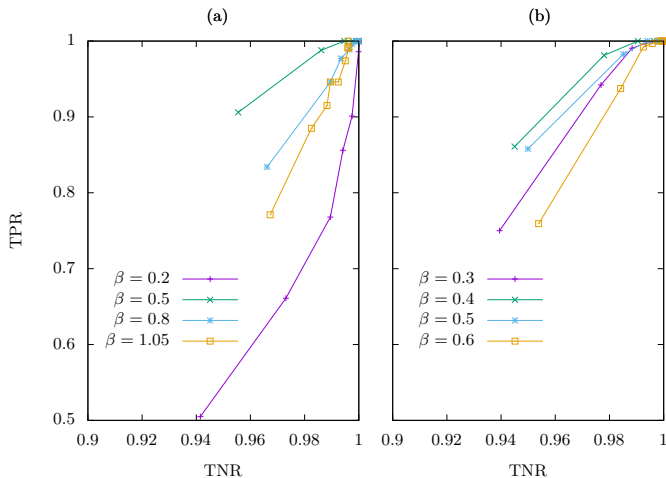


FIG. 4. ROC values for the (a) random graph case considered in Fig. (2), (b) diamond lattice case considered in Fig. (3). In both cases, each point corresponds to a different value of M , for $M = 100, 200, 300, 400, 500, 1000, 2000, 4000, 8000$, and they approach (1,1) as M increases.

the reconstructed graphs using this criterion is studied in detail in Fig. 4, where we plot the ROC parameters TNR and TPR for different temperatures, for the cases

discussed above.

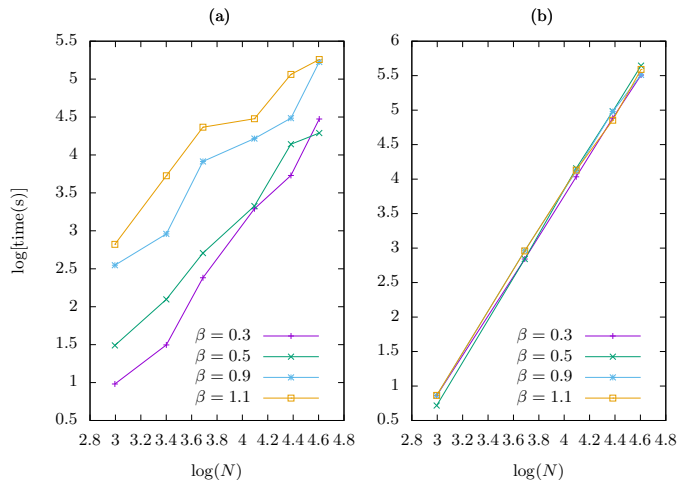


FIG. 5. Comparison between the log of the execution time (seconds) needed to find a solution as function of N and β for a RR spin glass with $c = 3$ using (a) PAMPL and (b) MPF. Each point corresponds to an average over 10 runs of the algorithm with $M = 4000$. $\beta_c = 0.615$.

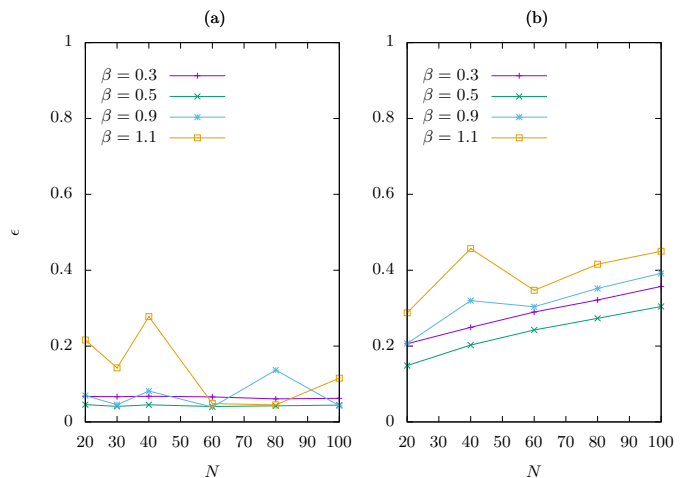


FIG. 6. Error ϵ , defined in eq. (9), as a function of N and β for a RR spin glass with $c = 3$ using (a) PAMPL and (b) MPF. Each point corresponds to an average over 10 runs of the algorithm at $M = 4000$.

We compare the performances of PAMPL with those of another fast inference method, namely the MPF [4]. MPF is as fast as a single maximization of the pseudo-likelihood and needs to be complemented with a threshold procedure. In Fig. 5-6 we analyze a RR graph with $c = 3$ with these two methods. MPF is expected to be $O(MN^2)$, as ours. While the tests show a more pronounced temperature dependence for PAMPL, we observe that the execution times are of the same order, both being very fast. Moreover, as the decimation algorithm improved the performances of methods based on

the maximization of the pseudo-likelihood, similarly the errors in the reconstructed graph made by PAMPL are 2-3 time smaller than those made by MPF. We also note that while the errors made by PAMPL does not depend on N , the error made by MPF do. More details on MPF and the case of a 2-D ferromagnetic lattice is discussed in details in Appendix D.

In summary we presented a new method, to reconstruct the hidden structure of Ising models based on the pseudo-likelihood and an activation procedure that includes recursively new parameters into a set whose elements are then optimized over. The method is exact in the limit of very large number of samples M and does not require setting ad-hoc extra parameters, apart from the choice of K which is mostly irrelevant. Performances of PAMPL are as good as or better than existing algorithms both based on PSL and other approaches, and the method is especially suitable to study inference problems with underlying sparse graphs.

This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation program (grant agreement No [694925]). S. Franz and J. Rocchi acknowledge the support of a grant from the Simons Foundation (No. 454941, Silvio Franz).

-
- [1] D H Ackley, G E Hinton, and T J Sejnowski. *Cognitive science*, 9(1):147–169, 1985.
- [2] Hilbert J. Kappen and Francisco de Borja Rodríguez. Efficient learning in boltzmann machines using linear response theory. *Neural Computation*, 10(5):1137–1156, 1998.
- [3] Toshiyuki Tanaka. Information geometry of mean-field approximation. *Neural Computation*, 12(8):1951–1968, 2000.
- [4] Jascha Sohl-Dickstein, Peter B Battaglino, and Michael R DeWeese. New method for parameter estimation in probabilistic models: minimum probability flow. *Physical review letters*, 107(22):220601, 2011.
- [5] S Cocco and R Monasson. *Physical review letters*, 106(9):090601, 2011.
- [6] E Aurell and M Ekeberg. *Physical review letters*, 108(9):090201, 2012.
- [7] F Ricci-Tersenghi. *Journal of Statistical Mechanics: Theory and Experiment*, 2012(08):P08015, 2012.
- [8] H C Nguyen and J Berg. *Physical review letters*, 109(5):050602, 2012.
- [9] S Cocco and R Monasson. *Journal of Statistical Physics*, 147(2):252–314, 2012.
- [10] J Raymond and F Ricci-Tersenghi. *Physical Review E*, 87(5):052111, 2013.
- [11] A Decelle and F Ricci-Tersenghi. *Physical review letters*, 112(7):070603, 2014.
- [12] A Y Lokhov, M Vuffray, S Misra, and M Chertkov. *Science advances*, 4(3):e1700791, 2018.
- [13] P Ravikumar, M J Wainwright, J D Lafferty, et al. *The Annals of Statistics*, 38(3):1287–1319, 2010.
- [14] Alessia Marruzzo, Payal Tyagi, Fabrizio Antenucci, Andrea Pagnani, and Luca Leuzzi. Inverse problem for multi-body interaction of nonlinear waves. *Scientific reports*, 7(1):3463, 2017.
- [15] Alessia Marruzzo, Payal Tyagi, Fabrizio Antenucci, Andrea Pagnani, and Luca Leuzzi. Improved pseudolikelihood regularization and decimation methods on nonlinearly interacting systems with continuous variables. *SciPost Physics*, 5(1):002, 2018.
- [16] Daniele Ancora and Luca Leuzzi. Learning direct and inverse transmission matrices. *arXiv preprint arXiv:1901.04816*, 2019.
- [17] Julian Besag. Efficiency of pseudolikelihood estimation for simple gaussian fields. *Biometrika*, pages 616–618, 1977.
- [18] G Schwarz et al. *The annals of statistics*, 6(2):461–464, 1978.
- [19] D C Liu and J Nocedal. *Mathematical programming*, 45(1-3):503–528, 1989.
- [20] As N Berker and S Ostlund. *Journal of Physics C: Solid State Physics*, 12(22):4961, 1979.
- [21] M Kaufman and R B Griffiths. *Physical Review B*, 24(1):496, 1981.

SUPPLEMENTAL MATERIAL

Appendix A: Pseudo-likelihood

We show that in the infinite sampling limit \mathcal{S} and \mathcal{L} are maximized by the same J . From the definition of \mathcal{S} in eq. (4), we obtain the property

$$\begin{aligned} \frac{\partial \mathcal{S}}{\partial J_{ij}} &= 2\beta \langle s_i s_j \rangle_{\text{Data}} - \beta \left\langle s_i \tanh \beta \sum_{m \neq j} J_{jm} s_m \right\rangle_{\text{Data}} \\ &\quad - \beta \left\langle s_j \tanh \beta \sum_{m \neq i} J_{im} s_m \right\rangle_{\text{Data}} \end{aligned} \quad (10)$$

and thus \mathcal{S} is maximum on the parameters J such that

$$\langle s_i s_j \rangle_{\text{Data}} = \left\langle s_i \tanh \beta \sum_{m \neq j} J_{jm} s_m \right\rangle_{\text{Data}}. \quad (11)$$

On the other hand, using the identity

$$\langle s_i s_j \rangle_{\text{Model}} = \left\langle s_i \tanh \beta \sum_{m \neq j} J_{jm} s_m \right\rangle_{\text{Model}} \quad (12)$$

in eq. (3), we notice that \mathcal{L} is maximum when

$$\langle s_i s_j \rangle_{\text{Data}} = \left\langle s_i \tanh \beta \sum_{m \neq j} J_{jm} s_m \right\rangle_{\text{Model}}, \quad (13)$$

and thus, since in the infinite sampling limit $\lim_{M \rightarrow \infty} \langle f(\underline{s}, J) \rangle_{\text{Data}} = \langle f(\underline{s}, J) \rangle_{\text{Model}}$, we observe that \mathcal{S} and \mathcal{L} are maximized by the same J .

Appendix B: Decimation algorithm

The idea of the decimation algorithm [11] is that starting from the complete graph, the full \mathcal{S} is maximized (maximization step) and the K couplings with the smallest values are set to zero (decimation step). The two steps are iterated until when no more couplings are present in the graph. In order to locate the stopping point, a new function is defined. Be \mathcal{S}_{max}^c the maximum of the pseudo-likelihood on the complete graph. When no couplings remain, the pseudo-likelihood is $-N \log 2$. The new function is given by $\mathcal{S}_t(x) = \mathcal{S} - [x\mathcal{S}_{max}^c - (1-x)N \log 2]$, where $x \in [0, 1]$ is the fraction of couplings, being 1 on the complete graph and 0 at the end of the decimation process. This function is 0 at the beginning and at the end of the process, by construction, and it is positive in the intermediate steps. The stopping point x^* is chosen by looking at the maximum of \mathcal{S}_t . The solutions found with this method are much better than those found with

other methods based on PSL. In Fig. 7 we show the behavior of this algorithm in the study of a 2-D ferromagnetic lattice with free boundary conditions with $N = 36$. In the inset we show the TPR and the TNR evolving with the iterations. The point where \mathcal{S}_t is maximum coincides with the point where ϵ is minimum. In Fig. 1 we analyze the same dataset with PAMPL and find a solution much faster because we start from the empty graph, rather than from the fully connected graph. Each step of the decimation algorithm is $O(MN^2)$, since it needs to optimize the PSL over the number of couplings that haven't been decimated yet. If the true graph is sparse we need to run the iterations for $O(N^2)$ times. In comparison, ours takes $O(MN^2)$ operations to provide a solution in sparse graphs, as explained in the following section and it is thus much faster.

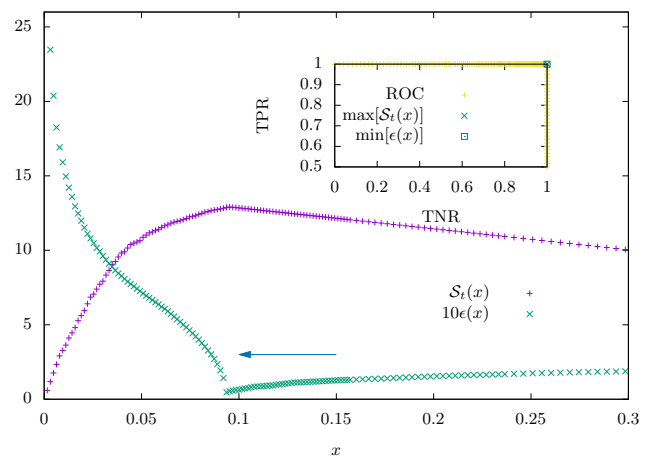


FIG. 7. 2-D ferromagnetic lattice with $N = 36$ and free boundary conditions, $M = 5000$, $\beta = 0.5$. Main Figure: \mathcal{S}_t and 10ϵ as a function of x for the decimation based method. x is 1 at the beginning and decreases as the iteration proceeds. Only the last part of the iteration is presented. The Arrow specifies the directions in which the iteration moves. Inset: ROC curve, with the TNR on the x-axis and the TPR on the y-axis.

Appendix C: Details of the implementation

Bayesian Information Criterion: In PAMPL we cannot use the tilted pseudo-likelihood to locate the stopping point because this would require the maximization of the pseudo-likelihood on the complete graph. As stated in the main text, the quantity that we observe during learning is thus the Bayesian Information Criterion, defined in eq. (8). Let us consider $P(\{\underline{s}^{(\mu)}\}_{\mu=1, \dots, M}) = \int dJ P(\{\underline{s}^{(\mu)}\}_{\mu=1, \dots, M} | J) P(J)$. Under the assumption of a flat prior $P(J)$, we expand $\mathcal{L}(J) = M^{-1} \log P(\{\underline{s}^{(\mu)}\}_{\mu=1, \dots, M} | J)$ to the second order around the parameters J^* for which the likelihood

$P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}|J)$ is maximum. A Gaussian integration leads to

$$P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}) = e^{M\mathcal{L}(J^*)} \left(\frac{2\pi}{I(J)M} \right)^{\frac{k}{2}} \quad (14)$$

where $I(J) = \mathcal{L}''(J^*) \sim O(1)$. Thus, we see that $P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}) \sim e^{M\mathcal{L}^* - k/2 \log M}$. In our analysis, we may replace \mathcal{L}^* with \mathcal{S}^* and obtain $P(\{\underline{s}^{(\mu)}\}_{\mu=1,\dots,M}) = e^{\text{BIC}/2}$, using eq. (8): the largest the BIC, the better the parameters J describe the observations.

Complexity: Updating and sorting the elements $\Delta\mathcal{S}_{ij}$, defined in eq. (4), requires a careful discussion. In fact, if we want to keep the iterative step $O(MN)$, we cannot afford $O(N^2)$ operations. The way we overcome this problem is explained below. At the initial time step we evaluate and sort all the gains in pseudo-likelihood $\Delta\mathcal{S}_{ij}$. After the sorting, we create the vector \underline{V} with the largest $O(N)$ elements. This vector is updated at each step and its size is kept to be $O(N)$. This makes the sorting less expensive. The cost of updating $O(N^2)$ element is alleviated by considering only $\Delta\mathcal{S}_{ij}$ whose nodes are involved in the activation of the couplings of the previous time step. More precisely, if coupling J_{ij} has been updated at time t , at $t+1$ we update $\Delta\mathcal{S}_{ik}$ and $\Delta\mathcal{S}_{jk}$ with $k=1,\dots,N$ and neglect the changes in the others. These operations cost $O(N)$. If some of these values happen to be larger than the mean value of the elements of \underline{V} , they are included in \underline{V} . Finally, all the elements $v_i < 0.01 \max\{v_i\}$ of \underline{V} are excluded from it. These wise precepts allows the size of \underline{V} to remain $O(N)$ and, thus, the ensemble of iterative steps to be $O(MN^2)$. The optimization of the couplings in the set \mathbb{J} is performed with a gradient ascent with a learning rate of 0.0001 until $\Delta\text{PSL}/\text{PSL} < 0.00001$, where ΔPSL is the difference between the PSL computed before and after the updating.

Derivatives: The expression of the first and second derivatives of \mathcal{S} with respect to J_{ij} in terms of the average over samples read

$$\frac{\partial\mathcal{S}}{\partial J_{ij}} = \frac{2\beta}{M} \sum_{\mu=1}^M s_i^{(\mu)} s_j^{(\mu)} \left[\frac{1}{1 + e^{2\beta s_i^{(\mu)} \sum_m J_{jm} s_m^{(\mu)}}} + \frac{1}{1 + e^{2\beta s_j^{(\mu)} \sum_m J_{im} s_m^{(\mu)}}} \right], \quad (15)$$

$$\frac{\partial^2\mathcal{S}}{\partial^2 J_{ij}} = -\frac{4\beta^2}{M} \sum_{\mu=1}^M \left[\frac{1}{2 + 2 \cosh \left[2\beta s_i^{(\mu)} \sum_m J_{jm} s_m^{(\mu)} \right]} + \frac{1}{2 + 2 \cosh \left[2\beta s_j^{(\mu)} \sum_m J_{im} s_m^{(\mu)} \right]} \right]. \quad (16)$$

It is easy to check that eq. (15) coincide with eq. (10).

Stopping point: The increase of the BIC after the correct stopping point is due to the fluctuations induced by the finite size sampling. This can be understood considering different datasets, each one made by M configurations. We use each dataset to extract the inferred graphs, and we observe the behavior of the BIC and the error. Then we compute the mean and the standard deviation of the two quantities and we observe that the minimum of the error is reached when the BIC reaches for the first time the value $\text{BIC}_{max} - \sigma_{\text{BIC}}$, being σ_{BIC} the standard deviation. In Fig. 8a we show the results for a 2-D lattice with $N=49$ and free boundary conditions. In Fig. 8b we show the results for a random regular graph with $N=100$ and mean connectivity equal to 4. In both cases we observe that the increase of the BIC after the correct stopping point (corresponding to the minimum of the error) is irrelevant. In order to locate the stopping point we thus adopt the criterion explained in the main text.

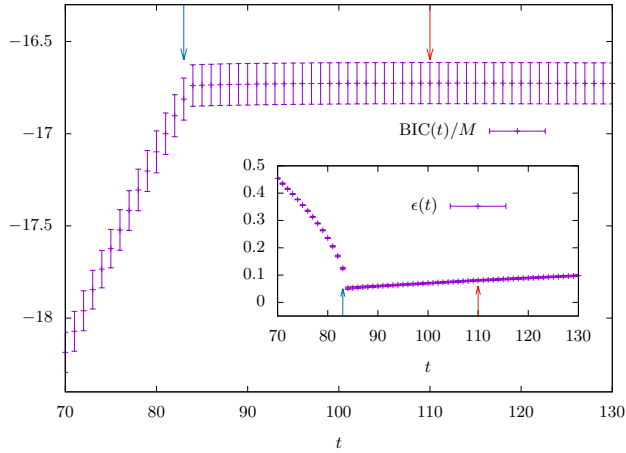
ROC curves: We consider a 2-D ferromagnetic lattice with periodic boundary condition in Fig. 10. In these figures we plot the TPR and the TNR for the inferred graph as a function of the (inverse of the) number of observed samples. We observe a weak dependence on the size of the system and a more severe one on the temperature. This is in line with the performances of specific algorithms for which it is possible to compute the scaling of the minimum number of samples for a perfect inference, where the dependence on N is logarithmic and that in β is exponential [12]. We notice that a TNR smaller than 1 means that the reconstructed graph contains couplings that are not present in the original graph, i.e. that our criterion does not detect correctly the stopping point and couplings keep being activated for some other step.

Appendix D: Minimum probability flow

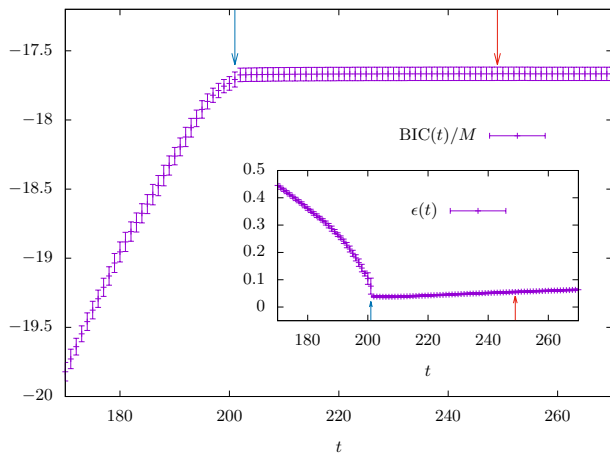
Algorithm: The Minimum Probability Flow (MPF) learning algorithm [4] is based on an hypothetical dynamics in the parameter space $\{J\}$, that we use to parametrize the probability distribution $P(\underline{s}|J) = \exp[-E_s(J)]/Z(J) = \exp\left(\sum_{i<j} J_{ij} s_i s_j\right)/Z(J)$. This dynamics starts from the data distribution and ends up in the point J that minimize the Kullback-Leibler D_{KL} divergence between the distribution of data and $P(\underline{s}|J)$. Using detailed balance, it is possible to define a transition matrix that allows the dynamics to relax to the chosen probability distribution,

$$\Gamma_{s,s'} = g_{s,s'} \exp\left[\frac{1}{2}(E_{s'} - E_s)\right] \quad (17)$$

with $g_{s,s'}$ being a sparse matrix with 1 between configurations differing by one-spin flip, and 0 elsewhere. The



(a)



(b)

FIG. 8. Average and standard deviation of the BIC over 10 runs of the activation-based algorithm. Blue and red arrows indicate the stopping point and the point at which the BIC is maximum. Main Figures: (a) 2-D spin glass lattice with $N = 49$ spins and 84 couplings with free boundary conditions. Each run of the algorithm is made on a subset of $M = 3000$ samples extracted from a dataset of 20000 samples at equilibrium at $\beta = 0.6$. (b) Random regular spin glass with $N = 100$ spins, $c = 4$, corresponding to 200 couplings. Each run of the algorithm is made on a subset of $M = 8000$ samples extracted from a dataset of 20000 samples at equilibrium at $\beta = 0.8$. Insets: Average and standard deviation of the error over the same $L = 10$ runs of the activation-based algorithm.

dynamics considered is thus

$$\partial_t p_s^{(t)} = \sum_{\underline{s}' \neq \underline{s}} \Gamma_{s,s'} p_{s'}^{(t)} - \sum_{\underline{s}' \neq \underline{s}} \Gamma_{s',s} p_s^{(t)}, \quad (18)$$

where $\Gamma_{s,s'}$ is the transition rate from configuration \underline{s}' to \underline{s} . This dynamics may take several time steps to converge to the desired distribution and it is not practical. Rather than waiting such a long time, MPF considers a small time $t = \epsilon$. In fact, among the trajectories that leads to

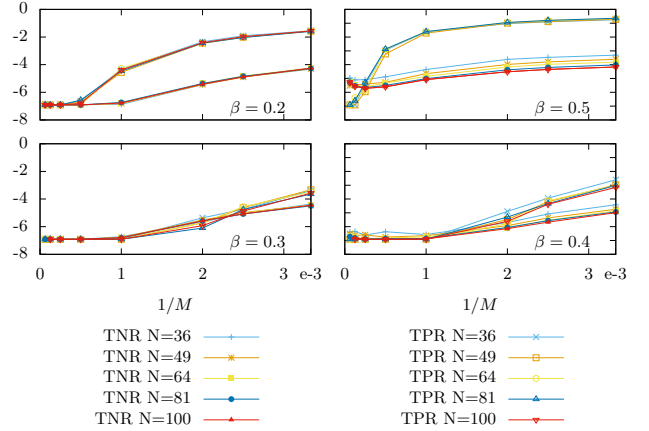


FIG. 9. $\log(1 - \text{TNR} + 0.001)$ and $\log(1 - \text{TPR} + 0.001)$ at different N for a 2D lattice with periodic boundary conditions at different β as a function of $1/M$. We notice that as M increases, both the TPR and TNR approach 1. Each point corresponds to the average of ~ 50 runs of the activation algorithm, stopped using the criterion defined in the text.

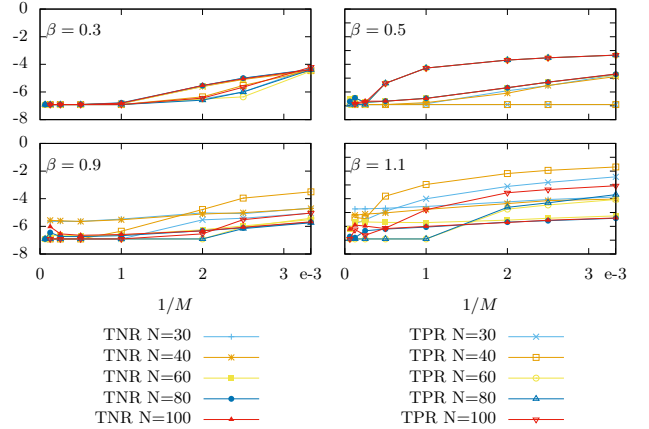


FIG. 10. $\log(1 - \text{TNR} + 0.001)$ and $\log(1 - \text{TPR} + 0.001)$ at different N for a RR graph with $c = 3$ at different β as a function of $1/M$. We notice that as M increases, both the TPR and TNR approach 1. Each point corresponds to the average of ~ 50 runs of the activation algorithm, stopped using the criterion defined in the text.

J^* , a special role is played by the one that points in the direction of J^* already in the early steps. In this limit, it is possible to show that

$$D_{KL}(p^{(0)}, p^{(t)}) = \frac{\epsilon}{M} \sum_{\underline{s}' \in \mathcal{D}} \sum_{\underline{s} \notin \mathcal{D}} \Gamma_{s,s'} \equiv K(J), \quad (19)$$

where \mathcal{D} denotes the dataset. Parameter estimation is provided by $J^* = \arg \min_J K(J)$. If the system is big enough, and the configurations of the dataset sampled independently, it is likely that configuration space is

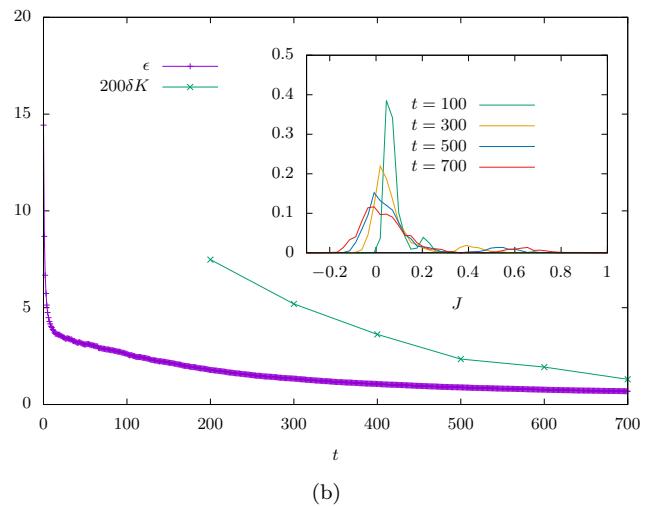
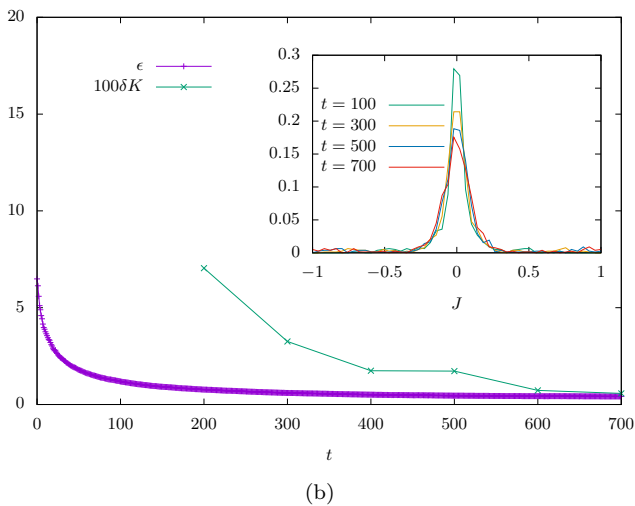
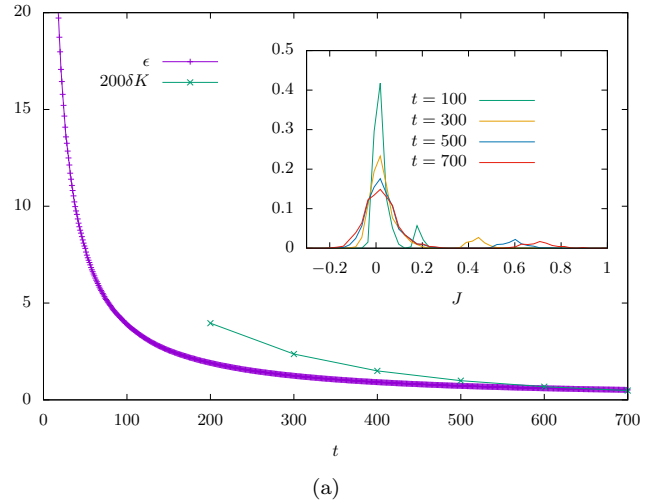
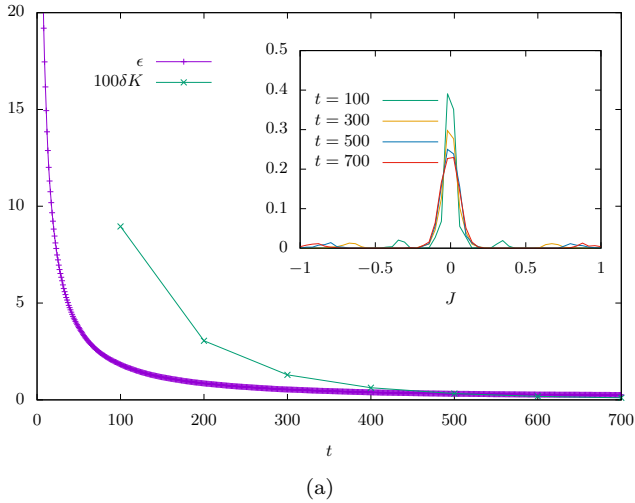


FIG. 11. ϵ and δK as a function of the iteration time t for the MPF learning on a RR graph with $c = 3$, $N = 40$ and $M = 4000$. (a) $\beta = 0.3$ (b) $\beta = 1.1$. The total execution time is roughly 30 seconds on a normal laptop, learning is stopped after 700 steps. Insets: histograms of the learned couplings J at different stages of learning. After the last time, δK is equal to 0.0011 (a) and 0.0057 (b). In both cases the majority of the true couplings have been found: setting a threshold at $|J| = 0.5$, TPR and TNR are respectively 1-1 (a) and 0.95-1 (b).

FIG. 12. ϵ and δK as a function of the iteration time t for the MPF learning on a 2D lattice with periodic boundary condition, $N = 49$ and $M = 4000$. (a) $\beta = 0.2$ (b) $\beta = 0.5$. The total execution time is roughly 55 seconds on a normal laptop, learning is stopped after 700 steps. Insets: histograms of the learned couplings J at different stages of learning. After the last time, δK is equal to 0.0065 (a) and 0.0025 (b). In both cases the majority of the true couplings have been found: setting a threshold at $|J| = 0.5$, TPR and TNR are respectively 1-1 (a) and 0.983-1 (b).

not sampled extensively and thus, for each configuration $s' \in \mathcal{D}$ of the dataset, all those that differ from it for a spin flip are not part of \mathcal{D} . The second sum in the definition of $K(J)$ is thus replaced by $\sum_{\underline{s}: g_{s, s'} = 1}$. This makes each step of the minimization process $O(MN)$. On the other hand, given that the algorithm optimizes over all the parameters J , we observe that the actual cost of each learning step is $O(MN^2)$. For the case under considera-

tion, $E_s = -\sum_{i < j} J_{ij} s_i s_j$ and

$$K(J) = \frac{\epsilon}{M} \sum_{\mu=1}^M \sum_{t=1}^N e^{-\beta \sum_{q \neq t} J_{tq} s_t^{(\mu)} s_q^{(\mu)}}. \quad (20)$$

$K(J)$ can be optimized either with a simple gradient descent and with a more sophisticated LBFGS [19] algorithm with similar results. Performances depend on the learning rate. In particular, when β is large, generally a smaller learning rate needs to be used. Moreover, a small mini-batch allows to find solutions in a smaller amount

of time. Mini-batch size should not be smaller than a few dozens in any case. A comparison between execution times in Figs 5-13 shows that MPF is as fast as PAMPL. In the paper we present results with $M = 4000$, 100 mini-batch and $\epsilon = 0.025$, that plays the role of a learning rate:

$$\frac{\partial K}{\partial J_{ij}} = -\frac{\epsilon\beta}{M} \sum_{\mu=1}^M s_i^{(\mu)} s_j^{(\mu)} \times \left[e^{-\beta \sum_{q \neq i} J_{iq} s_i^{(\mu)} s_q^{(\mu)}} + e^{-\beta \sum_{q \neq j} J_{jq} s_j^{(\mu)} s_q^{(\mu)}} \right]. \quad (21)$$

Moreover, we observe that MPF provides an alternative method to infer couplings that is as fast as a single maximization of pseudo-likelihood on the complete graph. This is more clear by a comparison with eq. (15): the updating rules used to maximize the MPF and the pseudo-likelihood \mathcal{S} are very similar, with the first one consisting in neglecting the two denominators $2 \cosh(\beta h_i)$ and $2 \cosh(\beta h_j)$ appearing in the second one, where they act as normalizations factors. As other methods based on the pseudo-likelihood, MPF needs to be complemented with a threshold procedure.

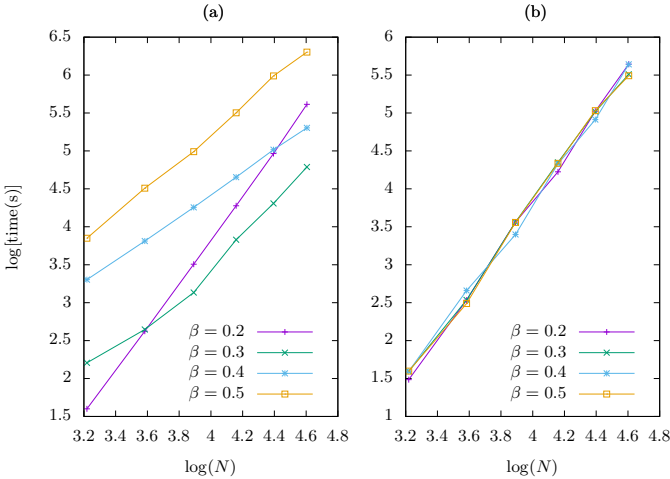


FIG. 13. Comparison between the log of the execution time (seconds) needed to find a solution as function of N and β for a 2D lattice with periodic boundary conditions using (a) PAMPL and (b) MPF. Each point corresponds to an average over 10 runs of the algorithm with $M = 4000$.

Stopping point

Finding a stopping point for this algorithm is not easy. In particular, even if the error with respect the original graph decreases quickly, the values of the couplings are still far from the actual ones and are refined only in later time steps. In particular, for real application cases where the actual topology is unknown, one should rely on other measures of convergence, like for instance $\delta K = \Delta K/K$,

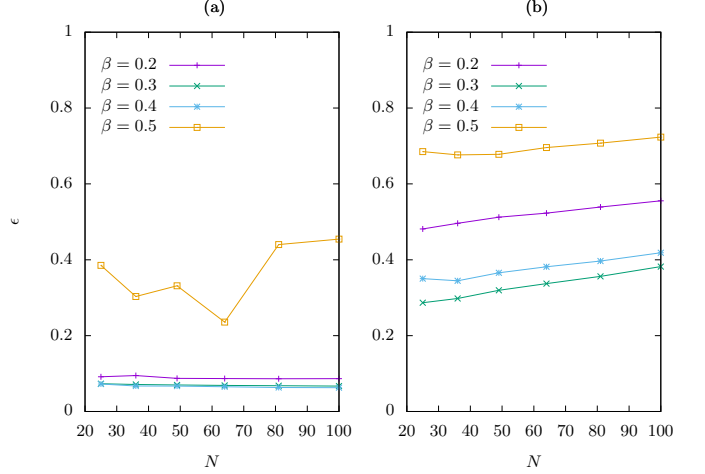


FIG. 14. Error ϵ , defined in eq. (9), as a function of N and β for a 2D ferromagnetic lattice with periodic boundaries using (a) PAMPL and (b) MPF. Each point corresponds to an average over 10 runs of the algorithm at $M = 4000$.

where ΔK is the difference on K computed every epoch running over the whole dataset. This quantity decreases during learning but our experiments do not provide a meaningful value where to stop the iteration. Unsurprisingly, parameters like mini-batch size, learning rate, stopping point are model, size and temperature dependent. In particular at large temperatures fewer iterations are needed to find satisfying results, as observed in the case of Fig. 11-12. As discussed above, MPF has to be complemented with a threshold procedure. In the main text we consider a RR spin glass with $c = 3$ while here we consider a 2D lattice with ferromagnetic interactions and periodic boundary conditions. In both cases, after ~ 700 iterations, setting a threshold at $|J|=0.5$, both the TPR and the TNR are very close to 1. On the other hand, corresponding errors are still large as seen in the main text in Fig. 6 for the RR case and here in Fig. 14 for the 2D lattice. This is not surprising. In fact, although being very versatile and fast, it performs a single optimization.