

PAPER: Interdisciplinary statistical mechanics

Nishimori meets Bethe: a spectral method for node classification in sparse weighted graphs

Lorenzo Dall'Amico^{1,*}, Romain Couillet^{1,2} and Nicolas Tremblay¹

¹ GIPSA-Lab, Université Grenoble Alpes, CNRS, Grenoble INP, France

² Laboratoire d'Informatique de Grenoble (LIG), Université Grenoble Alpes, France

E-mail: lorenzo.dall-amico@gipsa-lab.fr, romain.couillet@gipsa-lab.fr and nicolas.tremblay@gipsa-lab.fr

Received 5 March 2021

Accepted for publication 16 August 2021

Published 24 September 2021



Online at stacks.iop.org/JSTAT/2021/093405
<https://doi.org/10.1088/1742-5468/ac21d3>

Abstract. This article unveils a new relation between the *Nishimori temperature* parametrizing a distribution P and the *Bethe* free energy on random Erdős–Rényi graphs with edge weights distributed according to P . Estimating the Nishimori temperature being a task of major importance in Bayesian inference problems, as a practical corollary of this new relation, a numerical method is proposed to accurately estimate the Nishimori temperature from the eigenvalues of the *Bethe Hessian* matrix of the weighted graph. The algorithm, in turn, is used to propose a new spectral method for node classification in weighted (possibly sparse) graphs. The superiority of the method over competing state-of-the-art approaches is demonstrated both through theoretical arguments and real-world data experiments.

Keywords: clustering techniques, inference of graphical models, machine learning, random matrix theory and extensions

*Author to whom any correspondence should be addressed.

Contents

1. Introduction	2
1.1. From statistical physics...	3
1.2. ... to Bayesian inference on weighted sparse networks	3
1.3. Our contribution: relating Nishimori to Bethe	4
2. Basic properties of the RBIM	5
2.1. Phase diagram	5
2.2. Relevant properties at the Nishimori temperature	8
3. A relation between β_N and the Bethe free energy	9
3.1. Preliminaries	9
3.2. Main result	10
3.2.1. Arguments in support of claim	12
3.2.1.1. Dense graphs	12
3.2.1.2. Sparse graphs	15
3.3. The relation between β_N and the Bethe-Hessian matrix	17
3.3.1. The Bethe free energy	17
3.3.2. Phase diagram	18
3.4. Estimation of β_N from H_β, J	19
4. Application to node classification	21
4.1. A generative model for node classification	21
4.2. The Nishimori temperature-based node classification algorithm	22
4.3. Relation to other spectral methods	23
4.3.1. The weighted Laplacian matrix	23
4.3.2. The naïve mean field approach	25
4.3.3. The ‘spin glass Bethe-Hessian’	25
4.4. Application to real data classification	26
5. Conclusion	29
Acknowledgments	30
Appendix A. An explicit expression for the matrix $F(g)$	30
Appendix B. Algorithm implementation	31
References	33

1. Introduction

1.1. From statistical physics...

The physics of disordered systems [1] and Bayesian inference for graph learning [2] have long been shown to be tied by a deep connection that has given rise to a host of efficient physics-inspired algorithms [3–5]. A particularly telling example where this relation stands out is the so-called *teacher–student* scenario, in which a set of observed random variables are the outcome of a generative model (the *teacher*) with some hidden parameters to be learned by the *student* [6].

As an instrumental example, we consider in this article the problem of statistical inference on a graph in which the random variable observed by the student is a weighted, undirected graph. Specifically, given a realization of an Erdős–Rényi graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ with vertex set \mathcal{V} and edge set \mathcal{E} ,³ a random *weighted adjacency matrix* $J \in \mathbb{R}^{n \times n}$, with $J_{ij} = J_{ji} \neq 0$ only if (ij) is an edge of \mathcal{G} , is observed by the *student* whose task is to infer some latent variable of the generative model of J . The non-null entries of the matrix J are independently generated by the teacher according to the law

$$P(x) = p_0(|x|)e^{\beta_N x}, \quad (1)$$

for an arbitrary non-negative function $p_0(\cdot)$ and some $\beta_N > 0$, which we from now on refer to as the *Nishimori temperature*⁴ [7]. The Nishimori temperature naturally appears in statistical physics in the *random bond Ising model* (RBIM), in which the vector $\mathbf{s} \in \{-1, 1\}^n$ is a random variable distributed according to the Boltzmann distribution

$$\mu(\mathbf{s}) = \frac{e^{-\beta \mathcal{H}_J(\mathbf{s})}}{Z_{J, \beta}}, \quad (2)$$

for some positive β , with $Z_{J, \beta}$ a normalization constant and $\mathcal{H}_J(\mathbf{s}) = -\mathbf{s}^T J \mathbf{s}$.

At $\beta = \beta_N$, i.e. when the temperature of the system coincides with the Nishimori temperature⁵, the exact expression of $\mathbb{E}[\langle \mathcal{H}_J(\mathbf{s}) \rangle_\beta]$ can be computed with elementary mathematical tools, where $\langle \cdot \rangle_\beta$ denotes the averaging over the Boltzmann distribution (2) while $\mathbb{E}[\cdot]$ is the averaging over the realizations of J distributed as (1). It has also been shown [6, 8] that the RBIM at the Nishimori temperature is either in the *ferromagnetic* configuration (in which $\langle s_i \rangle_\beta > 0$ for all i) or in the *paramagnetic* configuration (for which $\langle s_i \rangle_\beta = 0$ for all i). In particular, the system is never in the *spin-glass phase* under which local order of \mathbf{s} appears despite there being no global magnetization. These relevant properties drew research attention to this particular temperature [9–11] since its first appearance in [7].

1.2. ... to Bayesian inference on weighted sparse networks

The importance of the *Nishimori temperature* in Bayesian inference was thoroughly discussed in [12], where the author exhibits a correspondence between the optimal Bayes

³The graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ is thus *fixed* and not a random variable.

⁴For the sake of precision, β_N behaves as an inverse temperature but, for simplicity, we will refer to it as a *temperature*.

⁵We underline here that, to be fully rigorous, β_N is not, by definition, a temperature, but rather a parameter of the generative model of J , i.e. a hidden parameter of the teacher's generative model.

inference problem (i.e. when the *student* knows exactly the generative model of the *teacher*) and the RBIM studied at β_N .

As a practical and telling example of modern concern of the importance of the Nishimori temperature in Bayesian statistics, we here consider as a common thread the problem of binary node classification on a graph. Specifically, let $\boldsymbol{\sigma} \in \{-1, 1\}^n$ be a label vector assigning each node to a ‘class’. Further assume that a matrix J is drawn from the distribution (1) and that the student has to infer the vector $\boldsymbol{\sigma}$ from the observation of the matrix \tilde{J} , defined by $\tilde{J}_{ij} = J_{ij}\sigma_i\sigma_j$. The matrix \tilde{J} has entries that, *in expectation*, are positive if nodes i and j have the same label and negative otherwise. As discussed extensively in section 4, this quite elementary model can in fact be used to study *correlation clustering* over the p -dimensional *feature* vectors $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^p$ of a dataset of size n [13], with concrete application to image, sound, or sentence classification [14]. In this example, the weights J_{ij} carried by the edges of \mathcal{G} represent some affinity metric between the features \mathbf{z}_i and \mathbf{z}_j associated with nodes i and j (in essence, the larger J_{ij} the closer \mathbf{z}_i and \mathbf{z}_j).

From a Bayesian perspective, inferring $\boldsymbol{\sigma}$ from \tilde{J} reduces to computing the marginals of the distribution

$$\mathbb{P}(\boldsymbol{\sigma}|\tilde{J}) = \frac{e^{-\beta_N \mathcal{H}_J(\boldsymbol{\sigma})}}{Z_{\tilde{J}, \beta_N}}. \quad (3)$$

This thus coincides with computing the *magnetizations* $\mathbf{m} = \langle \boldsymbol{\sigma} \rangle_{\beta_N}$ of an RBIM on the graph \tilde{J} at the Nishimori temperature. However, assuming that the observing *student* knows the value of β_N is often unrealistic (in effect, the student only sees \tilde{J}) and earlier works have resorted to studying the problem of *mismatched inference* (i.e. inference performed when the *student* uses a different parameter than the one assumed by the *teacher*) [6].

1.3. Our contribution: relating Nishimori to Bethe

Our main result consists in going beyond mismatched inference by providing an efficient estimate to the Nishimori temperature. To this end, we first draw an explicit relation between the Nishimori temperature and the smallest eigenvalue of the Hessian matrix of the Bethe free energy associated to the RBIM, when set at the paramagnetic point $\mathbf{m} = \mathbf{0}$ (this Hessian matrix is the so-called *Bethe-Hessian matrix* [15]); this relation holds under the previously introduced setting, so in particular for a student observation matrix J supported over a (possibly sparse) Erdős–Rényi graph \mathcal{G} . Besides, we observe and argue that, although the Bethe approximation is particularly adapted to sparse (locally tree-like) graphs \mathcal{G} , the *Nishimori–Bethe relation* holds for any degree of sparsity (that is, even when \mathcal{G} does not behave locally as a tree).

The main consequences of the Nishimori–Bethe relation, and our main contributions, consist in

- (a) The design of a new efficient spectral algorithm which estimates the Nishimori temperature with asymptotically perfect accuracy (as $n \rightarrow \infty$); the algorithm is based

on an iterative fast search of a well-parametrized Bethe-Hessian matrix exhibiting a smallest amplitude eigenvalue close to zero;

- (b) A new spectral algorithm to approximately solve the Bayesian node classification inference problem of equation (3), which outperforms commonly deployed state-of-the-art alternatives. We in particular claim that this spectral algorithm is capable of performing non trivial inference as soon as the Bayesian optimal solution can;
- (c) Although we claim that these algorithms are still valid under dense graphs \mathcal{G} , they are specifically adapted to the sparse regime where $|\mathcal{V}| \sim |\mathcal{E}|$; this practically allows for small computational and memory storage costs when applied to the classification of the nodes of possibly large graphs; we specifically support this fact by a concrete application to the classification of 40 000 high resolution images using our proposed sparse but extremely efficient spectral algorithm.

The remainder of the article is structured as follows. Section 2 introduces the RBIM together with some basic properties of the Nishimori temperature. These serve as the support for section 3, which provides our main results: the Nishimori–Bethe relation, the aforementioned new algorithms to estimate the Nishimori temperature, and how it provides an approximate (but still accurate) solution to the Bayesian inference problem. To corroborate the claims made in this section, section 4 applies the results to a concrete node classification problem involving realistic images produced by generative adversarial networks [16]. Section 5 closes the article laying out some limitations and possible directions of improvement of the present analysis.

A Julia implementation of our proposed algorithm as well as the codes used to produce the results of this article is available at github.com/lorenzodallamico/NishimoriBetheHessian.

Notation: vectors are denoted in bold face. The notation $\mathbf{1}_n$ indicates the all ones vector of size n . Scalar and matrices are in standard font, with matrices denoted by capital letters. The notation ‘ \circ ’ indicates the Hadamard entry-wise product between two matrices of same size. The notation ∂i indicates the neighbourhood of node i on the graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$: $\partial i = \{j \in \mathcal{V} : (ij) \in \mathcal{E}\}$.

2. Basic properties of the RBIM

In this section we provide the basic language and properties necessary to define the Nishimori temperature. The results presented in this section do not all have a direct application to inference problems, which are discussed later in section 4.

2.1. Phase diagram

Consider a realization of an Erdős–Rényi graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ with expected average degree c . We will denote \mathcal{V} the set of the n nodes of the graph and \mathcal{E} the set of its edges. We further let $J \in \mathbb{R}^{n \times n}$ be a weighted adjacency matrix on \mathcal{G} and distributed according to the following generative model.

Definition 1 (Generative model of \mathbf{J}). For all edges $(ij) \in \mathcal{G}$ with $i > j$, the J_{ij} are generated independently (with $J_{ij} = J_{ji}$), for some $\beta_N > 0$, referred to as *Nishimori temperature*, according to

$$\begin{aligned} \forall (ij) \in \mathcal{E}, \quad i < j, \quad J_{ij} \stackrel{\text{i.i.d.}}{\sim} P \\ P(x) = p_0(|x|)e^{\beta_N x}, \end{aligned} \tag{4}$$

where $p_0(\cdot)$ is an arbitrary non-negative function satisfying the normalization condition $\int_{-\infty}^{\infty} dx p_0(|x|)e^{\beta_N x} = 1$. If $(ij) \notin \mathcal{E}$, then $J_{ij} = 0$.

Given a realization of J and a vector $\mathbf{s} \in \{-1, 1\}^n$, we define the *Hamiltonian* $\mathcal{H}_J(\mathbf{s})$ of the RBIM as

$$\mathcal{H}_J(\mathbf{s}) = - \sum_{(ij) \in \mathcal{E}} J_{ij} s_i s_j = -\mathbf{s}^T J \mathbf{s}. \tag{5}$$

Note that, from definition 1, the Nishimori temperature is defined *independently* of \mathcal{G} , while the dependence of $p_0(\cdot)$ on β_N is relegated to its normalization constant. Two examples of distributions that fall under this definition are the $\pm J$ model

$$P(x) = p\delta(x - J_0) + (1 - p)\delta(x + J_0), \quad \text{for } p \in [1/2, 1], \quad J_0 \in \mathbb{R}^+$$

that can be rewritten as

$$P(x) = \frac{e^{\beta_N x}}{2\text{ch}(\beta_N J_0)}, \quad \text{with } \beta_N = \frac{1}{2J_0} \log \frac{p}{1 - p},$$

and the Edwards–Andersons model [17]

$$P(x) = \frac{1}{\sqrt{2\pi\nu^2}} \exp \left\{ -\frac{(x - J_0)^2}{2\nu^2} \right\}, \quad \text{for } J_0, \nu \in \mathbb{R}^+$$

for which

$$\begin{aligned} p_0(|x|) &= \frac{1}{\sqrt{2\pi\nu^2}} \exp \left\{ -\left(\frac{x^2 + J_0^2}{2\nu^2} \right) \right\} \\ \beta_N &= \frac{J_0}{\nu^2}. \end{aligned}$$

Given a matrix J drawn from the generative model of definition 1, equation (5) and a temperature $\beta \in \mathbb{R}^+$, we now let $\mathbf{s} \in \{-1, 1\}^n$ be a random vector, drawn from the Boltzmann distribution

$$\mu(\mathbf{s}) = \frac{e^{-\beta\mathcal{H}_J(\mathbf{s})}}{Z_{J,\beta}}, \tag{6}$$

where $Z_{J,\beta}$ is the normalization constant. Averaging over the distribution (6) will be denoted with $\langle \cdot \rangle_\beta$.

Let us now consider the phase diagram, depicted in figure 1, of the model described by equations (5) and (6) and definition 1. First consider the role played by the two

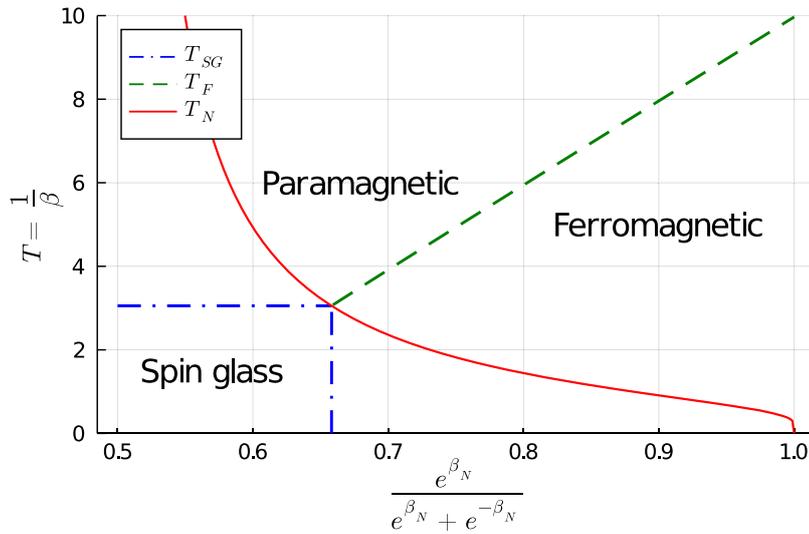


Figure 1. Phase diagram of the RBIM for $J_{ij} \in \{-1, 1\}$. The x axis goes from $\frac{1}{2}$ for $\beta_N = 0$ to 1 for $\beta_N \rightarrow \infty$. The y axis represents T , the inverse of β . The dashed green line is the inverse of β_F , the dash dotted blue line is the inverse of β_{SG} and the solid red line is the inverse of β_N .

parameters β and β_N . For increasing values of β_N , there is a larger probability for each edge J_{ij} to carry a positive weight and the minimum of $\mathcal{H}_J(\mathbf{s})$ is achieved for $\mathbf{s} = \mathbf{1}_n$. For small values of β_N , instead, multiple local minima appear. Concerning β , instead, for small values, the Boltzmann distribution (equation (6)) tends toward a uniform distribution, while, for large values, the configurations with small energy $\mathcal{H}_J(\mathbf{s})$ have a larger probability. Consequently, for large β and β_N the average configuration of \mathbf{s} tends to align toward $\mathbf{1}_n$: this corresponds to the *ferromagnetic* configuration. Conversely, for small values of β_N , several edges carry a negative weight, introducing *frustration* in the system that is found in the *spin-glass* phase, for which local order of the spins may be observed ($\frac{1}{n} \sum_i \langle s_i \rangle_\beta^2 \neq 0$), but globally the magnetization is null ($\frac{1}{n} \sum_i \langle s_i \rangle_\beta = 0$). Finally, at large values of β , the system is in the *paramagnetic* phase, for which the spins are randomly aligned and the magnetization is zero.

In the particular case where \mathcal{G} is an Erdős–Rényi random graph, with expected average degree equal to c , the cavity method [18] predicts the position of the transitions between the three phases: the *paramagnetic–ferromagnetic* transition occurs at $\beta = \beta_F$ and the *paramagnetic–spin glass* transition occurs at $\beta = \beta_{SG}$, also known as the de Almeida–Thouless transition [19]. The values of β_F, β_{SG} are given as the solutions of the following equations [6]:

$$c \cdot \mathbb{E}[\text{th}(\beta_F J_{ij})] := 1 \tag{7}$$

$$c \cdot \mathbb{E}[\text{th}^2(\beta_{SG} J_{ij})] := 1, \tag{8}$$

where we recall that $\mathbb{E}[\cdot]$ denotes averaging over the distribution (4). Figure 1 precisely depicts the phase diagram for the $\pm J$ model. A qualitatively similar diagram can be obtained for different distributions that follows the definition of equation (1) [7]. Given

these premises, we now discuss some relevant properties valid on the Nishimori line, i.e. when $\beta = \beta_N$.

2.2. Relevant properties at the Nishimori temperature

First of all, let us introduce the *quenched internal energy density*, defined as $u(\beta) := \frac{1}{n} \mathbb{E}[\langle \mathcal{H}_J(\mathbf{s}) \rangle_\beta]$, where we recall that $\langle \cdot \rangle_\beta$ denotes an average taken over the Boltzmann distribution (6) and $\mathbb{E}[\cdot]$ is the average over the distribution of equation (1). It was shown in [7] that $u(\beta_N)$ takes a particularly simple expression:

$$u(\beta_N) = \frac{1}{n} \mathbb{E}[\langle \mathcal{H}_J(\mathbf{s}) \rangle_{\beta_N}] = -\frac{1}{n} \sum_{(ij) \in \mathcal{E}} \mathbb{E}[J_{ij} \langle s_i s_j \rangle_{\beta_N}] = -\frac{1}{n} \sum_{(ij) \in \mathcal{E}} \mathbb{E}[J_{ij} \operatorname{th}(\beta_N J_{ij})].$$

The first two equalities are true by definition. The elegance of the result of [7] lies in the last relation that identifies—inside the expectation $\mathbb{E}[\cdot]$ —the term $\langle s_i s_j \rangle_{\beta_N}$ with $\operatorname{th}(\beta_N J_{ij})$. We will show in section 3.3 that, for sufficiently small β , the system is in the paramagnetic phase $\langle \mathbf{s} \rangle_\beta = 0$ and, under the Bethe approximation, the relation $\langle s_i s_j \rangle_\beta = \operatorname{th}(\beta J_{ij})$ is verified for any underlying β_N . This informally introduces a relation between the Bethe free energy at the paramagnetic point and the Nishimori temperature, which is at the center of claim 1.

We introduce the following property of the probability distribution of equation (1). This relation will be of fundamental use in the following and, in passing, allows us to rewrite $u(\beta_N)$ as in [7].

Property 1. *Let $f(x)$ be an arbitrary odd function. Then*

$$\mathbb{E}[f(x) \cdot \operatorname{th}(\beta_N x)] = \mathbb{E}[f(x)]. \tag{9}$$

The proof is easily obtained by straightforward calculation. As a consequence of property 1, the quenched internal energy density at the Nishimori temperature takes the simple expression:

$$u(\beta_N) = -\frac{1}{n} \sum_{(ij) \in \mathcal{E}} \mathbb{E}[J_{ij}] = -\frac{\bar{d}}{2} \cdot \mathbb{E}[J_{ij}],$$

where \bar{d} denotes the average node degree in the graph \mathcal{G} .

Secondly, we recall a well celebrated property of the Nishimori temperature, which states the absence of *replica symmetry breaking* on the Nishimori line [6, 8] or, equivalently, that the RBIM at the Nishimori temperature is never in the spin glass phase. This result can be visually understood in figure 1 by noticing that the Nishimori temperature is either in the paramagnetic or ferromagnetic phase. Moreover, exploiting property 1 and the definitions of β_F, β_{SG} in equations (7) and (8), on an Erdős–Rényi graph one finds that $\beta_{SG} = \beta_N \Leftrightarrow \beta_F = \beta_N$. Consequently, there exists a tricritical point where $\beta_F = \beta_{SG} = \beta_N$.

Recalling the connection with statistical inference problems, such as inferring σ in equation (3), first note that β_N is the Bayes optimal inference temperature in the sense that there exists no other β that can asymptotically achieve better inference

performance and, therefore, if inference cannot be performed at $\beta = \beta_N$, then it is theoretically infeasible. This occurs when the marginals of equation (3) asymptotically give equal probabilities for each σ_i to take either values ± 1 . In terms of the phase diagram, this corresponds to being in the *paramagnetic* phase, so that $\beta_N < \beta_F$. In order for non-trivial reconstruction to be possible, the condition $\beta_N < \beta_{SG} < \beta_F$ must be imposed [20]. When the condition is met, the system is in the *informative* configuration in which each *spin* gets oriented toward its planted value σ_i . This being said, replacing (or effectively, erroneously estimating) β_N by $\beta \neq \beta_N$ in equation (3), it may occur that, even though inference is theoretically possible (as $\beta_N < \beta_{SG} < \beta_F$), the estimated labels $\hat{\sigma}$ for the mismatched β are not aligned with the ground truth σ . This never happens at $\beta = \beta_N$ for which inference is achieved *as soon as theoretically possible*.

With this short introduction on the Nishimori temperature at hand, in the next section we present our main result which relates β_N to the spectrum of the non-backtracking and Bethe-Hessian matrices of the underlying graph \mathcal{G} .

3. A relation between β_N and the Bethe free energy

This section introduces our main theoretical result, which draws a connection between the Nishimori temperature and the *variational free energy under the Bethe approximation*, computed at the *paramagnetic* point $\langle \mathbf{s} \rangle_\beta := \mathbf{m} = 0$. To this end, section 3.1 introduces two fundamental matrices, namely the *non-backtracking* and the *Bethe-Hessian* matrices of the graph \mathcal{G} and recalls the known connections between the spectra of these two matrices. Section 3.2 then introduces our main result, precisely consisting in (i) a claim on the location of the eigenvalues of the *non-backtracking* matrix and, as a result of the claim, (ii) an explicit relation between the underlying Nishimori temperature and a specific eigenvalue of the *non-backtracking* matrix. We further provide both theoretical arguments and numerical simulations in support of the result. As a corollary of the identities listed in sections 3.1 and 3.2, we finally obtain an explicit relation between the *smallest eigenvalue of the Bethe-Hessian matrix* and the Nishimori temperature. Section 3.3 relates this central link to the phase diagram of figure 1, in passing connecting the results to the expression of the Bethe free energy. Based on these findings, section 3.4 elaborates an algorithm to estimate β_N , which finds significant importance in statistical inference problems.

3.1. Preliminaries

Let us first introduce the weighted *non-backtracking* matrix of any arbitrary graph \mathcal{G} .

Definition 2 (Weighted non-backtracking matrix). Given a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ and a function $f: \mathcal{E} \rightarrow \mathbb{R}$, so that $\forall e \in \mathcal{E}$, $f(e) = \omega_e$ is the weight corresponding to the edge e , the weighted non backtracking matrix $B \in \mathbb{R}^{2|\mathcal{E}| \times 2|\mathcal{E}|}$ is defined on the set of directed edges of \mathcal{G} as

$$B_{(ij),(k\ell)} = \delta_{jk}(1 - \delta_{i\ell})\omega_{k\ell}. \quad (10)$$

The *non-backtracking* matrix plays an important role in inference and graph mining problems [21–27] and naturally comes into play from the linearization of the *belief propagation* (BP) (or *cavity*) equations [18] for the RBIM. These equations are particularly adapted to dealing with locally *tree-like* structured graphs (such as sparse Erdős–Rényi graphs).

The eigenvalues of the matrix B are strongly related to the eigenvalues of the *Bethe-Hessian* matrix.

Definition 3 (Bethe-Hessian matrix). Given a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, a function $f: \mathcal{E} \rightarrow \mathbb{R}$ so that $\forall e \in \mathcal{E}, f(e) = \omega_e$ and a parameter $x \in \mathbb{C} \setminus \{\pm\omega_{ij}\}_{(ij) \in \mathcal{E}}$, the Bethe-Hessian matrix $H(x) \in \mathbb{C}^{n \times n}$ is defined as

$$H_{ij}(x) = \left(1 + \sum_{k \in \partial i} \frac{\omega_{ik}^2}{x^2 - \omega_{ik}^2} \right) \delta_{ij} - \frac{x \omega_{ij}}{x^2 - \omega_{ij}^2}. \tag{11}$$

Since \mathcal{G} is an undirected graph, $H(x)$ is symmetric but not Hermitian, unless $x \in \mathbb{R}$. The relation between the spectra of the matrices B and $H(x)$ is given by the Watanabe–Fukumizu formula [28, 29].

Property 2 (Watanabe–Fukumizu). Let $H(x)$ and B be defined as per (10) and (11) on the same graph \mathcal{G} and for the same weighting function f . Further let $x \in \mathbb{C} \setminus \{\pm\omega_{ij}\}_{(ij) \in \mathcal{E}}$. Then,

$$\det [xI_{2|\mathcal{E}|} - B] = \det [H(x)] \prod_{(ij) \in \mathcal{E}} (x^2 - \omega_{ij}^2), \tag{12}$$

so that, for all x in the spectrum of B , $\det[H(x)] = 0$.

With this preliminary information, we now proceed to the formulation of our main result which first consists in a conjecture on the spectrum of B , and which we then relate to the spectrum of $H(x)$ through property 2. Choosing $f(e) = \text{th}(\beta J_e)$ in the definition of B , where the weights J_e are distributed according to equation (4), we finally unfold the relation between the spectra of B , $H(x)$ and the Nishimori temperature.

3.2. Main result

We now proceed to studying the spectrum of the matrix B in the case where \mathcal{G} is an Erdős–Rényi graph and its weights ω_e (equation (10)) are drawn i.i.d. satisfying $|\omega_e| < 1$ with $\mathbb{E}[\omega] > 0$ sufficiently large. The interest of this setting in relation to the RBIM and the Nishimori temperature is to consider $\omega_e = \text{th}(\beta J_e)$ for $\beta_N > \beta_{SG}$ and J as per definition 1. In this particular case, one of the eigenvalues of B —and, as a consequence of property 2, a corresponding (more easily estimated) eigenvalue of the Bethe-Hessian matrix—has a direct relation with β_N .

The matrix B is not symmetric, hence its eigenvalues are not necessarily real. Since B is real though, the non-real eigenvalues come in complex-conjugate pairs. When weights are assigned independently at random in the interval $(-1, 1)$, we observe, in agreement with the theoretical results obtained on the spectrum of B [30–33], that in the $n \rightarrow \infty$ limit, the non-real eigenvalues of B are bounded by a circle on the complex plane and

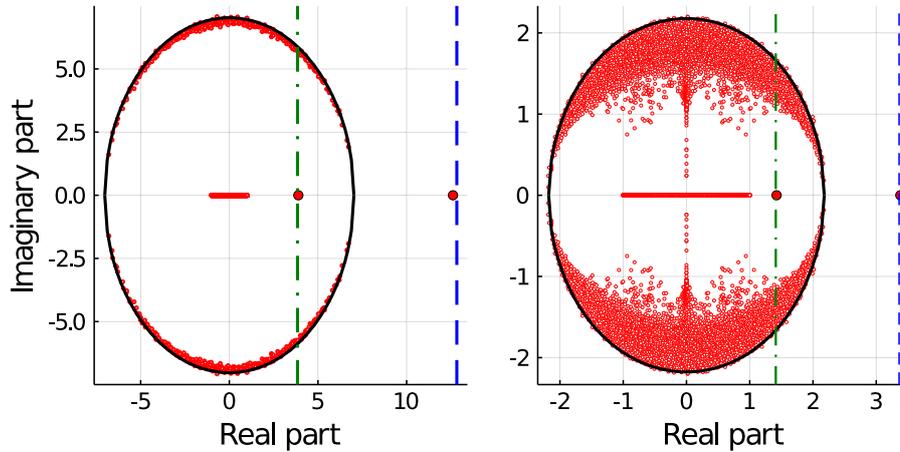


Figure 2. Spectrum of the matrix B in the complex plane. The entries J_{ij} are generated independently according to $\mathcal{N}(J_0, \nu^2)$. The weights appearing in equation (10) are defined as $\omega_{ij} = \text{th}(\beta J_{ij})$. (Left) Dense regime: $n = 250$, $c = 2 \log^2(n)$, $J_0 = 1$, $\nu = 4$, $\beta = 1$. (Right) Sparse regime: $n = 3000$, $c = 5$, $J_0 = 1$, $\nu = 1$, $\beta = 10$. For both plots, the dashed blue line corresponds to $c\mathbb{E}[\text{th}(\beta J)]$, the dash-dotted green line to $\mathbb{E}[\text{th}^2(\beta J)]/\mathbb{E}[\text{th}(\beta J)]$, while the black continuous line is the circle in the complex plane centred at the origin and of radius $\sqrt{c\mathbb{E}[\text{th}^2(\beta J)]}$.

are separated by a vanishing distance from one another. These eigenvalues form the *bulk* of the spectrum of B (see figure 2). There further exists one *real* eigenvalue which is instead *isolated*, i.e. it is found at a macroscopic (not decreasing with n) distance from all other eigenvalues. This eigenvalue has a modulus greater than the radius of the bulk: its existence and position are known and have been thoroughly investigated [32, 33]. There however exists another real isolated eigenvalue with modulus *smaller* than the radius of the bulk, the existence and importance of which were first evidenced in [34] in the case of unweighted graphs with a community structure. After [34], a similar phenomenon has also been observed in [35] in the context of phase retrieval, relating the Hessian of the TAP free energy and the Bayes optimal inference temperature. This isolated eigenvalue inside the bulk of B received less theoretical attention and it is the main object of our central result.

Claim 1. Let $\mathcal{G}(\mathcal{V}, \mathcal{E})$ be a realization of an Erdős–Rényi random graph with n nodes ($n \rightarrow \infty$) and expected average degree c . For each undirected edge $(ij) \in \mathcal{E}$ a weight $\omega_{ij} = \omega_{ji} \in (-1, 1)$ is assigned independently at random. Further assume that $\mathbb{E}[\omega_{ij}^2]/\mathbb{E}[\omega_{ij}] \geq 1$ and $\mathbb{E}[\omega_{ij}^2]/\mathbb{E}^2[\omega_{ij}] < c$. Then, the spectrum of B , with high probability, can be described as follows:

- There exist only two real eigenvalues in the spectrum of B with modulus greater or equal to one:

$$\lambda_1 = c\mathbb{E}[\omega] + o(c), \quad \lambda_{-1} = \frac{\mathbb{E}[\omega^2]}{\mathbb{E}[\omega]} + o(1). \tag{13}$$

The eigenvalue λ_1 is the largest in modulus in the spectrum of B ;

- All eigenvalues with non-zero imaginary part have a modulus bounded by $R = \sqrt{c\mathbb{E}[\omega^2]} + o(\sqrt{c})$.

Note that claim 1 does not make the assumption that $c \rightarrow \infty$ as $n \rightarrow \infty$, nor that $c = O_n(1)$. Extensive simulations indeed concur in suggesting that the claim holds in both dense and sparse graph regimes. The claim is thus stated for *any* average degree, so long that the underlying graph is of the Erdős–Rényi type. In detail, the technical condition $\mathbb{E}[\omega_{ij}^2]/\mathbb{E}^2[\omega_{ij}] < c$ is set to enforce that the leading eigenvalue λ_1 is greater than the radius of the bulk spectrum (hence that it is isolated) and that λ_{-1} is smaller than the radius of the bulk: a transition occurs at $\mathbb{E}[\omega_{ij}^2]/\mathbb{E}^2[\omega_{ij}] = c$ where both eigenvalues coincide: $\lambda_1 = \lambda_{-1}$. This inequality condition will thus ensure, when it comes to statistical inference, that non-trivial $\sigma \in \{\pm 1\}^n$ configurations can be theoretically recovered (i.e. that the inference problem is feasible). As a practical support to claim 1, figure 2 displays the spectrum of the matrix B in both moderately dense ($c \sim \log^2(n)$) and sparse ($c = O_n(1)$) regimes.

The fundamental corollary of claim 1 is that, in the case of present interest where $\omega_e = \text{th}(\beta J_e)$, from equation (13), the *inner* eigenvalue λ_{-1} of B is equal to

$$\lambda_{-1} = \frac{\mathbb{E}[\text{th}^2(\beta J)]}{\mathbb{E}[\text{th}(\beta J)]} + o(1).$$

Exploiting property 1, it follows immediately that, at $\beta = \beta_N$,

$$\lambda_{-1} \Big|_{\beta=\beta_N} = 1 + o(1).$$

Tuning the value of β until $\lambda_{-1} = 1$ thus provides *a method to estimate* β_N . The question on how to efficiently exploit this essential remark from an algorithmic standpoint will be further discussed in section 3.4.

Before pushing further our main line of deductions, we first introduce some arguments in support of claim 1, which we provide first in the dense and then in the sparse regimes. These are ‘arguments’ in the sense that they lack of full mathematical rigor and do not provide a formal proof of claim 1. Specifically, for the dense regime, we adopt a perturbative approach in which we heuristically show that the eigenvalues of B with modulus greater than one are close to the eigenvalues of the (easy to study) matrix M_0 appearing in equation (17). In the sparse regime, instead, we note that the position of the largest isolated eigenvalue and the radius of the bulk of B obtained in the dense case match the rigorous results of [32] proved for the sparse regime. With the support of extensive numerical simulations, we conjecture that the same result obtained in the dense regime to describe the *inner* isolated eigenvalue holds in the sparse regime as well.

3.2.1. Arguments in support of claim 1.

3.2.1.1. Dense graphs We first consider a dense graph regime, i.e. when the average degree c goes to infinity faster than $\log(n)$. This argument is inspired from the proof provided in [33] for unweighted dense graphs with a community structure. The proof of [33] can be straightforwardly adapted to the *binary* case in which $f(e) \in \{\pm\omega\}$, but does not unfold so directly for generic functions f .

The main advantage of the dense regime follows from the fact that the degree distribution of \mathcal{G} is almost regular and the Erdős–Rényi graph is close to a c -regular graph [36], the analysis of which is easier to handle. This makes it possible to relate the eigenvalues of B to those of $W \in \mathbb{R}^{n \times n}$, defined as $W_{ij} = \omega_{ij}$ if $(ij) \in \mathcal{E}$ and zero otherwise. The idea is to create a sequence of matrices $M(\mathbf{g}) \in \mathbb{R}^{2n \times 2n}$ (one for each eigenvector \mathbf{g} of B), in the spirit of a proof proposed by Bass of the celebrated Ihara–Bass formula [37], and to show that all the eigenvalues of $M(\mathbf{g})$ can be approximated, in the large n limit, by the eigenvalues of a common matrix M_0 independent of \mathbf{g} , so long that \mathbf{g} is an eigenvector corresponding to an eigenvalue λ of B for which $|\lambda| \geq 1$. It is the precise study of the spectrum of the limiting M_0 which induces the results of claim 1.

More specifically, let $\mathbf{g} \in \mathbb{C}^{2|\mathcal{E}|}$ be an eigenvector of B with eigenvalue λ , satisfying $|\lambda| \geq 1$ and let $\boldsymbol{\omega} \in \mathbb{R}^{2|\mathcal{E}|}$ be the vector containing the weights of the non-zero entries of the matrix B (and recall that $\omega_{ij} = \omega_{ji}$). Then define the vectors $\boldsymbol{\psi}(\mathbf{g}), \tilde{\boldsymbol{\psi}}(\mathbf{g}) \in \mathbb{C}^n$ as

$$\psi_i(\mathbf{g}) = \sum_{j \in \partial i} \omega_{ij} g_{ij}, \quad \tilde{\psi}_i(\mathbf{g}) = \sum_{j \in \partial i} \omega_{ij}^2 g_{ji} \tag{14}$$

and $F(\mathbf{g}) \in \mathbb{C}^{2n \times 2n}$ be any matrix satisfying

$$[F(\mathbf{g})\boldsymbol{\psi}(\mathbf{g})]_i = \sum_{j \in \partial i} \omega_{ij}^3 g_{ij}. \tag{15}$$

We now wish to relate the quantities $\boldsymbol{\psi}(\mathbf{g}), \tilde{\boldsymbol{\psi}}(\mathbf{g}), F(\mathbf{g})$ to the eigenvalues of B . In particular,

$$\begin{aligned} \lambda \psi_i(\mathbf{g}) &= \psi_i(B\mathbf{g}) = \sum_{j \in \partial i} \omega_{ij} \sum_{(kl)} \delta_{jk} (1 - \delta_{il}) \omega_{kl} g_{kl} = \sum_{j \in \partial i} \omega_{ij} \left[\sum_{l \in \partial j} \omega_{jl} g_{jl} - \omega_{ji} g_{ji} \right] \\ &= [W\boldsymbol{\psi}(\mathbf{g})]_i - \tilde{\psi}_i(\mathbf{g}) \end{aligned}$$

and, similarly,

$$\begin{aligned} \lambda \tilde{\psi}_i(\mathbf{g}) &= \tilde{\psi}_i(B\mathbf{g}) = \sum_{j \in \partial i} \omega_{ij}^2 \sum_{(kl)} \delta_{ik} (1 - \delta_{jl}) \omega_{kl} g_{kl} = \sum_{j \in \partial i} \omega_{ij}^2 \left[\sum_{l \in \partial i} \omega_{il} g_{il} - \omega_{ij} g_{ij} \right] \\ &= [D_W\boldsymbol{\psi}(\mathbf{g})]_i - [F(\mathbf{g})\boldsymbol{\psi}(\mathbf{g})]_i, \end{aligned}$$

where $D_W \in \mathbb{R}^{n \times n}$ is the diagonal matrix with $[D_W]_{ii} = \sum_{j \in \partial i} \omega_{ij}^2$. Thus, the eigenvalue λ is also an eigenvalue of the matrix

$$M(\mathbf{g}) = \begin{pmatrix} W & -I_n \\ D_W - F(\mathbf{g}) & 0 \end{pmatrix}. \tag{16}$$

The main difficulty of the analysis is of course introduced by the matrix $F(\mathbf{g})$ which is different for each eigenvector of B associated to $|\lambda| > 1$. In the *binary* case in which $W_{ij} \in \{\pm\omega\}$ for all $(ij) \in \mathcal{E}$, this term simplifies: combining equations (14) and (15), we get $F(\mathbf{g})\boldsymbol{\psi} = \omega^2\boldsymbol{\psi}$ and $F(\mathbf{g})$ thus simplifies for all \mathbf{g} into $F(\mathbf{g}) = \omega^2 I_n$; this allows

for a straightforward adaptation of the proof of [33]. The non-binary case is, however, more involved, but the term $(D_W - F(\mathbf{g}))\psi$ is still dominated by the action of D_W :

$$\left| \frac{[F(\mathbf{g})\psi(\mathbf{g})]_i}{\psi_i(\mathbf{g})} \right| = \left| \frac{\sum_{j \in \partial_i} \omega_{ij}^3 g_{ij}}{\sum_{j \in \partial_i} \omega_{ij} g_{ij}} \right| = o(c).$$

For the last equality, we exploited the fact that ω_{ij} and ω_{ij}^3 are both bounded in $(-1, 1)$ and have the same sign: this step is reasonable but non-rigorous, the main theoretical difficulty arising from the dependence between ω_{ij} and g_{ij} . Consequently, $F(\mathbf{g})$ can be regarded as a small perturbation of D_W . Further exploiting the concentration of the degrees, one may thus write

$$\|(D_W - F(\mathbf{g})) - c\mathbb{E}[\omega^2]I_n\| = o(c).$$

The eigenvalues of $M(\mathbf{g})$ can therefore be approximated by those of the matrix

$$M_0 = \begin{pmatrix} W & -I_n \\ c\mathbb{E}[\omega^2]I_n & 0 \end{pmatrix}. \tag{17}$$

The spectrum of M_0 is trivially related to the spectrum of W . Letting $\{\mu_i\}_{i=1, \dots, n}$ be the eigenvalues of W and $\{\lambda_{0,i}\}_{i=\pm 1, \dots, \pm n}$ those of M_0 , by the block determinant formula (section 5 of [38]), it comes that

$$\lambda_{0,\pm i} = \frac{\mu_i \pm \sqrt{\mu_i^2 - 4c\mathbb{E}[\omega^2]}}{2}. \tag{18}$$

In particular, it unfolds that

$$\begin{aligned} \mu_i^2 \geq 4c\mathbb{E}[\omega^2] &\implies \lambda_{0,-i} = \frac{c\mathbb{E}[\omega^2]}{\lambda_{0,i}} \equiv \frac{R^2}{\lambda_{0,i}} \\ \mu_i^2 < 4c\mathbb{E}[\omega^2] &\implies |\lambda_{0,\pm i}| = \sqrt{c\mathbb{E}[\omega^2]} \equiv R. \end{aligned}$$

Applying successively Wigner’s semi-circle theorem [39] and Bauer–Fike’s theorem [40], we thus have that

$$\mu_1 = c\mathbb{E}[\omega] + o(\sqrt{c}), \quad |\mu_i|_{|i| \geq 2} \leq \sqrt{c\mathbb{E}[\omega^2]} + o(\sqrt{c}). \tag{19}$$

Combining equations (18) and (19), along with the fact that the eigenvalues $\{\lambda_{0,\pm i}\}_{i=1, \dots, n}$ are a close approximation to the eigenvalues of B with modulus greater than one, we obtain the formulation of claim 1. Figure 3 compares the spectra of the matrices $M(\mathbf{g})$ and M_0 , which should be themselves compared to the left display in figure 2. Appendix A provides the explicit expressions of the matrix $F(\mathbf{g})$ used in figure 3.

This technical argument provides important intuitions on the spectrum of B : (i) the leading eigenvalue of B (the largest in modulus) is determined by the expectation of the entries of W ; (ii) the radius of the bulk of B is determined by the expectation of the squared entries of W ; (iii) all eigenvalues of B with modulus greater than one come in pairs (see equation (18)): they are complex conjugates if their imaginary part is non-zero

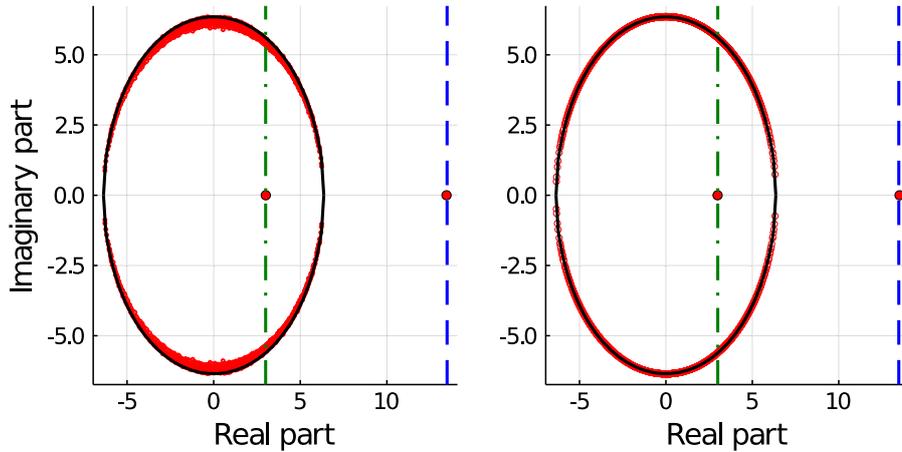


Figure 3. (Left) Spectrum of the matrices $M(\mathbf{g})$ defined in equation (16) with \mathbf{g} , one of the eigenvectors of B attached to a complex eigenvalue. (Right) Spectrum of M_0 , defined in equation (17). The graph considered is the same for the two matrices, with $n = 1500$, $c = \log^2(n)$. The matrix $W = \text{th}(\beta J)$, with $\beta = 1$ and the entries J_{ij} are i.i.d. normal variable with $J_0 = 1$ and $\nu = 3$. The blue dotted line is the position of $c\mathbb{E}[\text{th}(\beta J)]$, the green dash-dotted line is the position of $\mathbb{E}[\text{th}^2(\beta J)]/\mathbb{E}[\text{th}(\beta J)]$, while the black solid line is the circle in the complex plane of radius $\sqrt{c\mathbb{E}[\text{th}^2(\beta J)]}$.

or *harmonic conjugate* if they are real. This last observation justifies the existence of a real isolated eigenvalue *inside* the bulk of B , the importance of which will be further discussed in section 3.3.

As a downside, the setting considered in this section is, somehow, too simplistic. The analysis allowed us to neglect the term $F(\mathbf{g})$, which plays the role of the *Onsager reaction term* [41] which does not appear in the naïve mean-field approximation but plays a crucial role in the Bethe approximation. The fact that $F(\mathbf{g})$ can be neglected thus indicates that the regime under consideration is somehow *too simple*, the spectral behaviour of B being fully determined by W .

Consequently, we next discuss the far more interesting *sparse* regime in which the *Onsager reaction term* plays a fundamental role and in which the Bethe approximation brings a decisive advantage over the naïve mean-field approximation. In the sparse regime, the structure of the spectrum of the matrix B is essentially preserved, as well as the fact that all its eigenvalues all come in real harmonic or complex conjugate pairs.

3.2.1.2. Sparse graphs The Bethe approximation is exact on trees [18] and asymptotically yields (in the large n limit) exact results on locally *tree-like* graphs. This is precisely the case of sparse Erdős–Rényi graphs, in which the average degree is of order $c = O_n(1)$. In this case, the spectrum of the matrix W is no longer formed by an isolated eigenvalue (the largest in modulus) and a *bulk* of eigenvalues close to each other that follow the semi-circle law, as it happens in the *dense* regime discussed in the previous paragraph. Here the eigenvalues of W are known to have an unbounded support (little else is in fact theoretically known about this spectrum).

The non-backtracking matrix B , instead, essentially preserves the same spectral structure as in the dense regime, in which the *bulk* eigenvalues are bounded by a circle in the complex plane, as shown in figure 2. This result was recently proved in [32] in which the authors showed that, under the assumptions of claim 1, the matrix B has an isolated eigenvalue equal to $\lambda_1 = c\mathbb{E}[\omega] + o_n(1)$, (recall that $c = O_n(1)$) while all other eigenvalues satisfy $|\lambda_{i \geq 2}| \leq \sqrt{c\mathbb{E}[\omega^2]} + o(1)$. The result of [32] however does not mention the existence of *inner* real eigenvalues in the spectrum of B and, to best of our knowledge, no mathematical tool has been developed yet to rigorously address this question in the sparse regime. Yet, the position of the leading eigenvalue of B and the radius of its bulk spectrum are the same as in the *dense* graph case. We then conjecture, supported by extensive simulations, that also the inner isolated eigenvalue has the same position as in the dense regime, given by the square radius of the bulk, divided by the leading eigenvalue of B .

We take the opportunity of the reference to [32] to generalize the central claim of the article to their richer context. This result is of independent interest, particularly for more structured graph models.

Remark 1 (Random sparse graphs with independent entries). Note that the result of [32] is given under more general hypotheses than those discussed here. Specifically, the authors of [32] consider a setting in which each edge of the graph \mathcal{G} is created independently at random with probability p_{ij} . The Erdős–Rényi graph falls into the particular case in which $P = \{p_{ij}\}_{i,j=1}^n = \frac{c}{n} \mathbf{1}_n \mathbf{1}_n^T$. The leading (real) eigenvalues of B are determined from the leading eigenvalues of $P \circ \mathbb{E}[W]$, and the bulk radius by the leading eigenvalue of $P \circ \mathbb{E}[W \circ W]$. Based on this result, we conjecture that the real eigenvalues of B come in harmonic pairs precisely determined by

$$\lambda_{\pm i} = \frac{\rho(P \circ \mathbb{E}[W \circ W])}{\gamma_i(P \circ \mathbb{E}[W])},$$

where $\rho(\cdot)$ indicates the largest eigenvalue in modulus, and $\gamma_i(P \circ \mathbb{E}[W])$ are the eigenvalues of $P \circ \mathbb{E}[W]$, greater than $\sqrt{\rho(P \circ \mathbb{E}[W \circ W])}$.

A particular case of this setting is the *degree-corrected stochastic block model* which reproduces a k -class structure on an unweighted graph. In this case, the matrix $\mathbb{E}[W]$ has a low rank factorization $\mathbb{E}[W] = \frac{1}{n} \sum_{i=1}^k \alpha_i \mathbf{u}_i \mathbf{u}_i^T$, with $\alpha_1 = c$. Furthermore $P = \boldsymbol{\theta} \boldsymbol{\theta}^T$, where $\boldsymbol{\theta}^T \mathbf{1}_n = n$ and $\frac{1}{n} \boldsymbol{\theta}^T \boldsymbol{\theta} := \Phi$. Then, in agreement with [30] and the conjecture of [34], the eigenvalues of B can be described as follows

$$\begin{aligned} \lambda_i &= \alpha_i \Phi + o_n(1) \quad \text{for } 1 \leq i \leq k \\ \lambda_{-i} &= \frac{c}{\alpha_i} + o_n(1) \quad \text{for } 1 \leq i \leq k \\ |\lambda_{i > k}| &\leq \sqrt{c\Phi} + o_n(1). \end{aligned}$$

Returning to the implications of claim 1 of immediate interest, recall that the claim makes it possible to relate the Nishimori temperature to the specific eigenvalue λ_{-1} of the matrix B . From a numerical standpoint though, λ_{-1} is not easily accessible for

two reasons: (i) the matrix B is non-symmetric and large, slowing down eigenvalue computations; (ii) since λ_{-1} is smaller in modulus than most of the complex eigenvalues of B , while not being the smallest in modulus (see figure 2), one needs to compute all the bulk eigenvalues of B in order to access λ_{-1} : this comes at an impractical computational cost of $O(cn^3)$ with state of the art methods (see, for example, [42]). We next show that, as a consequence of claim 1 and property 1, the (symmetric) Bethe-Hessian matrix $H(x)$ (11) can be efficiently used to estimate β_N in the RBIM with a computational cost scaling as $O(nc)$.

3.3. The relation between β_N and the Bethe-Hessian matrix

This section elaborates on our final relation between the Bethe-Hessian matrix and the Nishimori temperature, as well as on how the respective spectra of the matrices $H(x)$ and B can be related to the phase diagram of figure 1.

3.3.1. The Bethe free energy. Let us first recall the basics of a variational approach, and specifically of the Bethe approximation. For $\mu(\mathbf{s})$ the Boltzmann distribution (6), the free energy $F_{J,\beta}$ and the variational free energy $\tilde{F}_{J,\beta}(\mathbf{q})$ (given for an arbitrary set of parameters \mathbf{q}), are defined through

$$F_{J,\beta} = \sum_{\mathbf{s}} \mu(\mathbf{s}) (\beta \mathcal{H}_J(\mathbf{s}) + \log \mu(\mathbf{s})) \tag{20}$$

$$\tilde{F}_{J,\beta}(\mathbf{q}) = \sum_{\mathbf{s}} p_{\mathbf{q}}(\mathbf{s}) (\beta \mathcal{H}_J(\mathbf{s}) + \log p_{\mathbf{q}}(\mathbf{s})). \tag{21}$$

The function $F_{J,\beta}$ is a moment generating function for the Boltzmann distribution of equation (2) but, in general, cannot be computed exactly. The variational free energy $\tilde{F}_{J,\beta}(\mathbf{q})$ represents a tractable approximation of $F_{J,\beta}$. From a straightforward calculation it can in particular be shown that $\tilde{F}_{J,\beta}(\mathbf{q}) - F_{J,\beta} = D_{\text{KL}}(\mu \| p_{\mathbf{q}}) \geq 0$, where $D_{\text{KL}}(\cdot \| \cdot)$ is the Kullback–Leibler divergence between two distributions. For a parametrized family of distributions $p_{\mathbf{q}}$, minimizing the variational free energy with respect to \mathbf{q} provides the Kullback–Liebler optimal approximation of $F_{J,\beta}$. The variational Bethe approximation considers a mean- and covariance-parametrized distribution $p_{\mathbf{q}} = p_{\mathbf{m},\chi}$ defined as

$$p_{\mathbf{m},\chi}(\mathbf{s}) = \prod_{(ij) \in \mathcal{E}} \frac{1 + m_i s_i + m_j s_j + \chi_{ij} s_i s_j}{4} \cdot \prod_{i=1}^n \left[\frac{1 + m_i s_i}{2} \right]^{1-d_i}, \tag{22}$$

where m_i and χ_{ij} are the average of s_i and $s_i s_j$ according to $p_{\mathbf{m},\chi}(\mathbf{s})$, respectively. Here d_i denotes the degree of node i ($d_i = |\{j : (ij) \in \mathcal{E}\}|$). The approximation turns out to be the exact factorization of $\mu(\mathbf{s})$ when \mathcal{G} is a tree, and is thus often claimed a good approximation of it in sparse, tree-like graphs.

A complete derivation of the Bethe-Hessian matrix from the Bethe free energy is proposed in [15]. It is instructive though to recall its main steps which allow one to relate the Bethe-Hessian matrix eigenvalues to the phase diagram of figure 1. From the expression of $\tilde{F}_{J,\beta}^{\text{Bethe}}(\mathbf{m}, \chi)$, obtained combining equations (21) and (22), one obtains that $\nabla_{\mathbf{m}} \tilde{F}_{J,\beta}^{\text{Bethe}}(\mathbf{m}, \chi)|_{\mathbf{m}=0} = 0$, i.e. the *paramagnetic* point is always an extremum of

the Bethe free energy. In order to study the stability of this solution, we consider the *Hessian* matrix of the variational free energy, computed at the paramagnetic point: the smallest eigenvalues of this matrix are associated to the *local* directions along which the paramagnetic solution may get unstable and non-trivial order in the spin configurations can be observed. This Hessian matrix explicitly reads:

$$\left. \frac{\partial^2 \tilde{F}_{J,\beta}^{\text{Bethe}}(\mathbf{m}, \boldsymbol{\chi})}{\partial m_i \partial m_j} \right|_{\mathbf{m}=\mathbf{0}} = \delta_{ij} \left(1 + \sum_{k \in \partial i} \frac{\chi_{ik}^2}{1 - \chi_{ik}^2} \right) - \frac{\chi_{ij}}{1 - \chi_{ij}^2}. \quad (23)$$

By further computing the gradient of $\tilde{F}_{J,\beta}^{\text{Bethe}}(\mathbf{m}, \boldsymbol{\chi})$ with respect to $\boldsymbol{\chi}$, one next obtains $\chi_{ij} = \text{th}(\beta J_{ij})$ (as already mentioned in the section 2 where we pointed that, in the paramagnetic phase, $\langle s_i s_j \rangle_\beta = \text{th}(\beta J_{ij})$ under the Bethe approximation). Setting $\omega_{ij} = \text{th}(\beta J_{ij})$, the matrix of equation (23) precisely corresponds to $H(1)$ defined in equation (11). We denote this matrix $H_{\beta,J}$, which explicitly reads:

$$(H_{\beta,J})_{ij} = \delta_{ij} \left(1 + \sum_{k \in \partial i} \frac{\text{th}^2(\beta J_{ik})}{1 - \text{th}^2(\beta J_{ik})} \right) - \frac{\text{th}(\beta J_{ij})}{1 - \text{th}^2(\beta J_{ij})}. \quad (24)$$

We may now relate the Bethe approximation to the phase diagram of figure 1.

3.3.2. Phase diagram. Let us move back to the system described by equations (5) and (6) and definition 1, first set at sufficiently high temperature (small β). In this case, for all β_N the system is in the *paramagnetic* phase, for which $\langle s_i \rangle_\beta = 0$. The paramagnetic solution $\mathbf{m} = \mathbf{0}$ is a minimum of $\tilde{F}_{J,\beta}^{\text{Bethe}}(\mathbf{m}, \boldsymbol{\chi})$, $H_{\beta,J}$ is positive definite.

Consider now β_N to be sufficiently large, so that the system undergoes to a *paramagnetic-ferromagnetic* phase transition (see figure 1). For $\beta = \beta_F$ defined as $c\mathbb{E}[\text{th}(\beta_F J)] = 1$, the leading eigenvalue of B is equal to 1 and one of the eigenvalues of $H_{\beta,J}$ is equal to zero. This eigenvalue is necessarily the smallest, since for $\beta < \beta_F$ all the eigenvalues are positive.

For small values of β_N , the system undergoes the *paramagnetic-spin glass* phase transition (see figure 1) at the temperature $\beta = \beta_{SG}$ defined so that $c\mathbb{E}[\text{th}^2(\beta_{SG} J)] = 1$. For this value of β , the radius of the bulk of the matrix B is equal to one and the bulk of $H_{\beta,J}$ is asymptotically close to zero.

Finally, further decreasing the temperature, at $\beta = \beta_N$ defined by $\mathbb{E}[\text{th}^2(\beta_N J)] = \mathbb{E}[\text{th}(\beta_N J)]$, the eigenvalue λ_{-1} is equal to one and the smallest eigenvalue of $H_{\beta,J}$ reaches zero for the second time. In figure 4, we show the spectra of the matrices B and $H_{\beta,J}$ at $\beta < \beta_F, \beta = \beta_F, \beta_{SG}, \beta_N$ that confirm the relation between the spectra of these two matrices and the phase diagram.

Having established in depth the relation between the matrices B and $H(x)$, and their relations to the phase diagram, we now show how one can efficiently estimate β_N , exploiting the smallest eigenvalue of $H_{\beta,J}$.

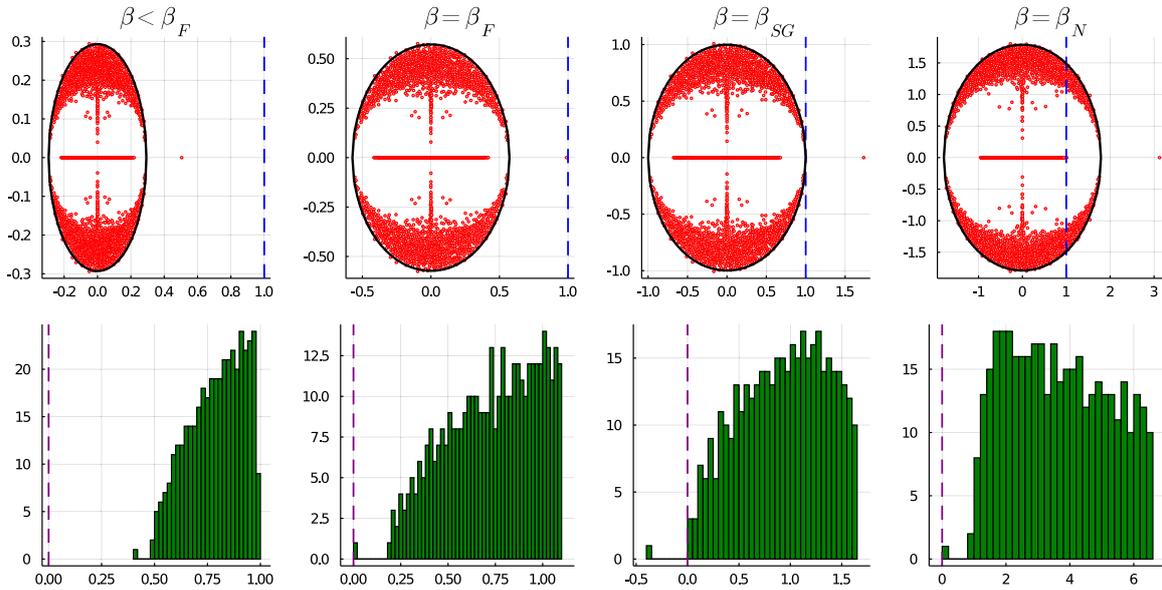


Figure 4. (First row) Spectrum of the matrix B in the complex plane for different values of β ; (second row) histogram of the eigenvalues of $H_{\beta,J}$ (zoomed in on the smallest eigenvalues) for different values of β . (First column) $\beta = 0.5\beta_F$, paramagnetic phase; (second column) $\beta = \beta_F$ paramagnetic–ferromagnetic transition; (third column) $\beta = \beta_{SG}$ paramagnetic–spin glass phase transition; (fourth column) $\beta = \beta_N$, Nishimori temperature. For all matrices, the same graph was used with $n = 1000$, $c = 10$. The weights of the edges are $\omega_{ij} = \text{th}(\beta J_{ij})$ for the different values of β just described. The J_{ij} are drawn independently from a Gaussian distribution with $J_0 = 1$ and $\nu = 1.5$. The blue lines in the first row is the vertical line at $x = 1$, while the purple line in the second row is the vertical line at $x = 0$.

3.4. Estimation of β_N from $H_{\beta,J}$

The present section provides a numerically efficient estimator $\hat{\beta}_N$ of the Nishimori temperature, first defined formally and then under the form of the output of a practical *algorithm*.

The proposed value of $\hat{\beta}_N$, estimate of the genuine Nishimori temperature β_N , reads

$$\hat{\beta}_N = \max_{\beta} \{ \beta : \gamma_{\min}(H_{\beta,J}) = 0 \}, \tag{25}$$

where $\gamma_{\min}(\cdot)$ indicates the smallest eigenvalue of a matrix. Under this definition, not only does $\hat{\beta}_N$ provides a consistent estimate of β_N for J distributed as definition 1, this being a consequence of claim 1, but it also provides the ‘best guess’ of an hypothetically corresponding β_N for matrices J which would follow a different distribution from the model of definition 1. Indeed, $\hat{\beta}_N$ has the advantage of always being defined, even for arbitrary matrices J , while having a clear interpretation for the class of matrices that fall under definition 1. This definition is particularly reminiscent of the algorithm proposed in [43] for community detection over sparse heterogeneous graphs, and which demonstrates a robust behaviour on applications to real-world graphs.

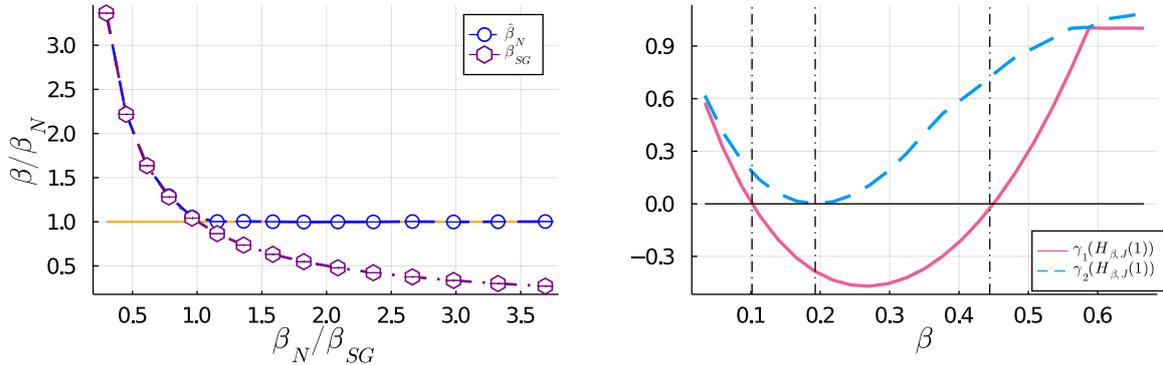


Figure 5. (Left) Computation of $\hat{\beta}_N$ for different values of β_N . The blue dots represent the ratio between $\hat{\beta}_N$, computed with algorithm 1 and the analytical value of β_N . The purple hexagons are the value of β_{SG}/β_N , while the orange line is at $y = 1$. For these plots, $n = 10\,000$ and $c = 5$. The weights of the non-zero entries of J are distributed i.i.d. according to $\mathcal{N}(J_0, \nu^2)$ for J_0 ranging from $J_0 = 0.5$ to $J_0 = 4$ and $\nu = 3.5$. Averages are taken over ten samples. (Right) Behaviour of the two smallest eigenvalues of $H_{\beta,J}$ as a function of β . The solid line indicates the smallest eigenvalue, while the dotted line is the second smallest. The vertical lines are set at $\beta_F < \beta_{SG} < \beta_N$. For this simulation, $n = 30\,000$ and $c = 10$. The weights of the matrix J are distributed i.i.d. according to a $\mathcal{N}(J_0, \nu^2)$ with $J_0 = 1$ and $\nu = 1.5$.

To best understand the rationale behind the definition of $\hat{\beta}_N$, first observe that $H_{\beta,J}$ is positive definite for all small values of β ($\lim_{\beta \rightarrow 0} H_{\beta,J} = I_n$). By increasing β , the smallest eigenvalue eventually hits zero a first time before turning negative: the zero-crossing occurs precisely at $\beta = \beta_F$. Then, continuing increasing β , at $\beta = \beta_{SG}$, the second smallest eigenvalue of $H_{\beta,J}$ is asymptotically equal to zero and $\gamma_{\min}(H_{\beta,J})$ is now negative. Finally, for $\beta \rightarrow \infty$, $H_{\beta,J}$ is again positive definite (the result can be easily obtained using Gershgorin’s circle theorem). Therefore, there must exist a value $\beta > \beta_{SG}$ for which $\gamma_{\min}(H_{\beta,J}) = 0$ for a second time. This second zero-crossing occurs precisely when $\beta = \hat{\beta}_N$. The right display of figure 5 visually explains this behaviour.

The basic idea of the proposed algorithm to compute $\hat{\beta}_N$ precisely consists in starting from $\beta = \beta_{SG}$ to then find the value of $\beta > \beta_{SG}$ for which $\gamma_{\min}(H_{\beta,J}) = 0$. Following this argument, we propose an iterative algorithm based on Courant–Fischer theorem to compute $\hat{\beta}_N$. The output of algorithm 1 is depicted in the left display of figure 5. Note in particular that, as long as $\beta_{SG} < \beta_N$, i.e. so long that $\mathbb{E}[\text{th}^2(\beta J_{ij})]/\mathbb{E}^2[\text{th}(\beta J_{ij})] < c$, the value of $\hat{\beta}_N$ is a good estimate of β_N . When the condition is instead not met, $\hat{\beta}_N$ simply coincides with β_{SG} . A more detailed analysis of algorithm 1 is provided in appendix B. The numerical advantage of exploiting the Bethe-Hessian matrix is decisive. First $H_{\beta,J}$ is symmetric and of size $n \times n$ regardless of the average node degree. Most importantly, the only eigenvalue of $H_{\beta,J}$ that needs be computed is the one of smallest amplitude, so that $\hat{\beta}_N$ can be estimated at an $O(nc)$ computational cost (using the Arnoldi method [42]).

While from a purely physics standpoint, claim 1 is an elegant theoretical relation between the Nishimori temperature and the Bethe-Hessian matrix, when it comes to

Algorithm 1. Compute $\hat{\beta}_N$.

input : Weighted adjacency matrix of a graph $J \in \mathbb{R}^{n \times n}$, precision error $\epsilon \in \mathbb{R}$;
output: Value of $\hat{\beta}_N \in \mathbb{R}^+$;
 Compute c , the average degree of the underlying unweighted graph: $c = \frac{1}{n} \sum_i \sum_j \mathbb{I}(J_{ij} \neq 0)$;
 Compute $\hat{\beta}_{SG}$ by solving $c \mathbb{E}[\text{th}^2(\hat{\beta}_{SG} J_{ij})] = 1$;
 Set $t = 1$ and $\beta_t \leftarrow \hat{\beta}_{SG}$;
 Initialize $\delta \leftarrow +\infty$;
while $\delta > \epsilon$ **do**
 Compute $H_{\beta_t, J}$ (equation (24)) ;
 Compute $\gamma_{\min, t}$, the smallest eigenvalue of $H_{\beta_t, J}$, as well as its associated eigenvector \mathbf{x}_t ;
 Define the function $f_t(\beta') = \mathbf{x}_t^T H_{\beta', J} \mathbf{x}_t$, for $\beta' \in \mathbb{R}^+$;
 Compute β_{t+1} by solving $f_t(\beta_{t+1}) = 0$;
 Update $\delta \leftarrow |\gamma_{\min, t}|$;
 Increment $t \leftarrow t + 1$;
return: β_{t-1}

machine learning applications, estimating β_N may have practical impact on algorithm performance. In particular, $\hat{\beta}_N$ may be used as an approximation of β_N when solving statistical inference on σ (for instance, via an optimal linearization of the Bayes optimal solution) in the absence of knowledge of the parameters in the generative model (3).

4. Application to node classification

This section discusses one of the immediate applications of the results introduced in the previous sections to the context of Bayesian statistical inference, and specifically to the problem of unsupervised node clustering on a graph. To this end, we first establish the relation between the Bayesian optimal inference and the Nishimori temperature, specific to the node classification problem; this then allows us to particularize algorithm 1 to this setting. Possibly most importantly, we conclude by commenting on how the considered model may be extrapolated to perform clustering on (possibly sparse) adjacency matrices of *real data* and relate our resulting proposed algorithm to commonly used competing spectral algorithms.

4.1. A generative model for node classification

Let \mathcal{G} be the realization of an Erdős–Rényi graph whose nodes are divided in two non-overlapping classes, labeled via the vector $\sigma \in \{-1, 1\}^n$. Associated to \mathcal{G} is a weighted adjacency matrix $\tilde{J} \in \mathbb{R}^{n \times n}$ with probability distribution:

$$\mathbb{P}(\tilde{J} | \sigma) = \prod_{(ij) \in \mathcal{E}} p_0(|\tilde{J}_{ij}|) e^{\beta_N \tilde{J}_{ij} \sigma_i \sigma_j}, \tag{26}$$

for an arbitrary non negative function $p_0(\cdot)$ and for some $\beta_N > 0$. According to this model, the edges connecting nodes in the same community are more likely to be positive,

while those connecting nodes in opposite communities are instead more likely to be negative. Given a realization of \tilde{J} , the task of the experimenter (who only has access to \tilde{J}) is to infer the vector σ . We can formulate this problem in terms of a Bayesian inference:

$$\mathbb{P}(\sigma|\tilde{J}) = \frac{\mathbb{P}(\tilde{J}|\sigma)\mathbb{P}(\sigma)}{\mathbb{P}(\tilde{J})} = \frac{1}{Z_{\tilde{J}}} \exp \left\{ \sum_{(ij) \in \mathcal{E}} \beta_N \tilde{J}_{ij} \sigma_i \sigma_j \right\}. \quad (27)$$

Computing the marginals of $\mathbb{P}(\sigma|\tilde{J})$ is equivalent to computing the average magnetization of an Ising model on \tilde{J} at the Nishimori temperature. However, the value of β_N cannot be easily inferred from \tilde{J} without knowing σ : one would indeed need to solve

$$\mathbb{E}[\text{th}(\beta_N \tilde{J}_{ij} \sigma_i \sigma_j)] = \mathbb{E}[\text{th}^2(\beta_N \tilde{J}_{ij} \sigma_i \sigma_j)].$$

To progress further, let us next introduce the matrix $J = \tilde{J} \circ \sigma \sigma^T$. Given the probability distribution of \tilde{J} (26), the matrix J is exactly defined as per definition 1. The key result to proceed consists in observing that the matrices $H_{\beta,J}$ and $H_{\beta,\tilde{J}}$ have the same eigenvalues and, up to a *gauge* transformation, the same eigenvectors. This then enables the use of algorithm 1 to estimate β_N directly from \tilde{J} . Let indeed $\mathbf{x} \in \mathbb{R}^n$ be an eigenvector of $H_{\beta,\tilde{J}}$ with eigenvalue λ and let \mathbf{y} have entries $y_i = x_i \sigma_i$. Then

$$\lambda y_i = \lambda x_i \sigma_i = \sigma_i \sum_j (H_{\beta,\tilde{J}})_{ij} x_j = \sigma_i \sum_j (H_{\beta,J})_{ij} \sigma_i \sigma_j x_j = (H_{\beta,J} \mathbf{y})_i$$

so that λ is an eigenvalue of $H_{\beta,J}$ with eigenvector \mathbf{y} . Consequently, the smallest eigenvalue of $H_{\beta_N,\tilde{J}}$ is asymptotically close to zero and algorithm 1 can be used to estimate β_N .

4.2. The Nishimori temperature-based node classification algorithm

For the purpose of node clustering though, the knowledge of β_N is a necessary prerequisite to obtain a precise estimate of the genuine node classes σ . We indeed show next that a powerful estimator of σ is obtained directly from the signs of the entries of the eigenvector \mathbf{x} of the Bethe-Hessian matrix $H_{\beta_N,\tilde{J}}$ (so β_N needs be known) associated to its smallest amplitude eigenvalue (which we now know is close to zero).

To this end, let us first consider \mathbf{y} , the eigenvector associated to the smallest eigenvalue of $H_{\beta_N,J}$. Denote with $A \in \{0, 1\}^{n \times n}$ the symmetric adjacency matrix of \mathcal{G} , defined by $A_{ij} = 1$ if $(ij) \in \mathcal{E}$, and $A_{ij} = 0$ otherwise, and let $D \in \mathbb{N}^{n \times n}$ be the diagonal degree matrix $D = \text{diag}(A \mathbf{1}_n)$. Then, applying property 1, one easily obtains that

$$\mathbb{E}[H_{\beta_N,J}] = I_n + \mathbb{E} \left[\frac{\text{th}(\beta_N J_{ij})}{1 - \text{th}^2(\beta_N J_{ij})} \right] (D - A). \quad (28)$$

From a straightforward calculation (see proposition 1 of [44]), the vector $\mathbf{1}_n$ is the eigenvector of $\mathbb{E}[H_{\beta_N,J}]$ associated to its eigenvalue of smallest amplitude. As a consequence, from the relation between \mathbf{x} and \mathbf{y} (or equivalently between \tilde{J} and J) in the previous section, the vector σ is the eigenvector of $\mathbb{E}[H_{\beta_N,\tilde{J}}]$ associated with its eigenvalue

Algorithm 2. The Nishimori–Bethe relation for node classification.

input: Weighted adjacency matrix of a graph $\tilde{J} \in \mathbb{R}^{n \times n}$, precision error $\epsilon \in \mathbb{R}$;
output: Value of $\hat{\beta}_N \in \mathbb{R}^+$, estimated label vector $\hat{\sigma} \in \{-1, 1\}^n$;
 Shift the non-zero \tilde{J}_{ij} as: $\tilde{J}_{ij} \leftarrow \tilde{J}_{ij} - \frac{1}{2|\mathcal{E}|} \mathbf{1}_n^T \tilde{J} \mathbf{1}_n$;
 Compute $\hat{\beta}_N \leftarrow \text{Compute_}\hat{\beta}_N$ (algorithm 1);
 Compute $H_{\hat{\beta}_N, \tilde{J}}$ (equation (24));
 Compute $\mathbf{x} \leftarrow$ the eigenvector associated to $\gamma_{\min}(H_{\hat{\beta}_N, \tilde{J}})$;
 Estimate $\hat{\sigma}$ as the output of two-class *k-means* on the entries of \mathbf{x} ;
return: $\beta_t, \hat{\sigma}$.

of smallest amplitude. Consequently, the eigenvector with zero eigenvalue of $H_{\hat{\beta}_N, \tilde{J}}$ is a close approximation⁶ of σ .

This conclusion immediately translates into algorithm 2, a numerical method to infer the genuine node classification σ . Further detail on the practical implementation of this algorithm are provided in appendix B.

Having established a ‘Nishimori-optimal’ version of the Bethe Hessian-based spectral clustering for node classification, the next section discusses the relation between the proposed algorithm and other commonly used kernel matrices in the spectral clustering literature.

4.3. Relation to other spectral methods

In the following, we use the *overlap*

$$\text{Overlap} = \left| 2 \left(\frac{1}{n} \sum_{i=1}^n \delta_{\sigma_i, \hat{\sigma}_i} - \frac{1}{2} \right) \right|. \tag{29}$$

as a measure of comparison of the inference performance of various node classification algorithms, where $\hat{\sigma}_i$ is the estimated label of node i . The overlap ranges from 0 (random assignment) to 1 (perfect assignment). Figure 6 compares the overlap achieved by algorithm 2 versus the naïve mean field approach, consisting in estimating the labels from the dominant eigenvector of \tilde{J} , and versus the popular legacy spectral clustering algorithm based on the weighted graph Laplacian matrix $L = \bar{D} - \tilde{J}$, where $\bar{D} = \text{diag}(|\tilde{J}| \mathbf{1}_n)$ [45].⁷ The figure browses several values of β_N (the larger β_N , the easier the detection problem) and of the average degree c . For $c = 3, 15$ the output of the asymptotically optimal BP algorithm is further shown, evidencing that algorithm 2 achieves an almost optimal performance. Due to its computational complexity, we chose not to run BP for $c = 50$, but we expect to observe a similar result to the one obtained for $c = 3, 15$. Before discussing the achieved results, let us first justify our comparison choice by recalling the rationale behind the Laplacian and naïve mean-field approaches.

⁶Rigorously, it is not so straightforward to move from $\mathbb{E}[H_{\beta_N, \tilde{J}}]$ to $H_{\beta_N, \tilde{J}}$. In [34], a similar setting is considered in which the eigenvector \mathbf{x} is studied in depth. The article argues that the relation indeed holds.

⁷Here $|\cdot|$ is the entry-wise absolute value.

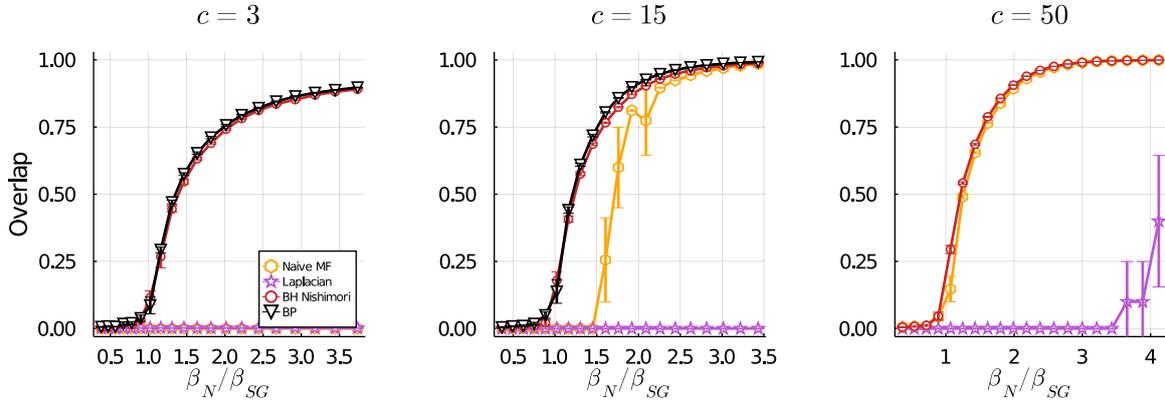


Figure 6. Overlap performance as a function of β_N/β_{SG} and three different values of the expected average degree, c . For $\beta_N < \beta_{SG}$ inference is asymptotically unfeasible. Two classes of equal size are considered and the entries of J are generated independently according to a Gaussian with mean $J_0\sigma_i\sigma_j$. In the examples, $n = 30\,000$ and n average is taken over 10 simulations.

4.3.1. The weighted Laplacian matrix. A very classical spectral clustering method in weighted graphs (dating back from the earliest works on the subject [44]) exploits the *weighted Laplacian* matrix $L = \bar{D} - \tilde{J}$, where $\bar{D} = \text{diag}(|\tilde{J}|\mathbf{1}_n)$. As shown in [44, 45], the eigenvector attached to the smallest eigenvalue of L provides a (discrete to continuous) relaxed solution of the NP-hard optimization *signed ratio-cut* graph clustering problem. The idea underlying the signed ratio-cut procedure consists in inferring the label assignments σ by maximizing the number of edges with positive weights connecting nodes in the same community, while minimizing the number of edges with negative weights connecting nodes in opposite communities: this is however a discrete optimization problem, a continuous relaxation of which coincides with a minimal eigenvector problem for L .

A particularly immediate and best understood scenario is the case of *signed graphs*, in which the entries of \tilde{J} assume values in ± 1 . For this class of graphs, an explicit relation between the matrices L and $H_{\beta_N, \tilde{J}}$ arises in the limit of trivial clustering, i.e. as $\beta_N \rightarrow \infty$. For signed graphs, a slightly different definition of $H_{\beta, \tilde{J}}$ than (24) is most appropriate:

$$H_{\beta, \tilde{J}}^{\text{signed}} = (1 - \text{th}^2(\beta))I_n + \text{th}^2(\beta)D - \text{th}(\beta)\tilde{J}. \quad (30)$$

It is straightforward to notice that the signed and unsigned versions of the Bethe-Hessian matrix share the same set of eigenvectors *on a signed graph* while their eigenvalues only differ by a multiplicative constant. One then immediately finds that $\lim_{\beta_N \rightarrow \infty} H_{\beta_N, \tilde{J}}^{\text{signed}} = L$. The signed Laplacian may then be seen as the *zero temperature limit* of the Bethe-Hessian matrix. From a Bayesian inference standpoint (27), $H_{\beta_N, \tilde{J}}^{\text{signed}}$ is a linear approximation of the exact inference problem, while L is only an approximation for the *maximum a posteriori* probability problem⁸. Far from the limit of trivial recovery, our proposed Bethe-Hessian matrix-based method is thus expected to accomplish better inference

⁸Taking the limit $\beta_N \rightarrow \infty$ in equation (27) is equivalent to looking for the maximum *a posteriori* solution.

performance when compared to the weighted Laplacian approach. This is indeed confirmed by figure 6, which evidences a striking performance gap between both methods. The reconstruction performance achieved through the matrix L is in particular severely compromised in the sparse regime in which \tilde{J} only has $O_n(n)$ non-zero entries.

4.3.2. The naïve mean field approach. The ‘Nishimori Bethe-Hessian’ matrix is built from the Bethe approximation of the Bayes optimal problem formulation. We now show that a similar approximation procedure could have been performed using a naïve mean field approximation instead. This leads to a different—much less efficient as we will see—spectral clustering algorithm. Recalling the procedure of section 3.3.1, we define the naïve mean field free energy from the probability distribution

$$p_{\mathbf{m}}(\mathbf{s}) = \prod_{i \in \mathcal{V}} \frac{1 + m_i s_i}{2}, \tag{31}$$

where m_i is the average of s_i over the distribution (31). The associated variational free energy reads

$$\tilde{F}_{\tilde{J}, \beta}^{\text{MF}}(\mathbf{m}) = - \sum_{(ij) \in \mathcal{E}} \beta \tilde{J}_{ij} m_i m_j + \sum_{i \in \mathcal{V}} \sum_{s_i} \frac{1 + m_i s_i}{2} \log \left(\frac{1 + m_i s_i}{2} \right).$$

Computing the gradient of $\tilde{F}_{\tilde{J}, \beta}^{\text{MF}}(\mathbf{m})$, one finds that, also in this case, the paramagnetic point $\mathbf{m} = \mathbf{0}$ is an extreme. Computing the Hessian of the free energy at the paramagnetic point leads instead to

$$H_{\beta, \tilde{J}}^{\text{MF}} = I_n - \beta \tilde{J}. \tag{32}$$

As a consequence, despite the presence of β in the formulation of $H_{\beta, \tilde{J}}^{\text{MF}}$, the eigenvectors of $H_{\beta, \tilde{J}}^{\text{MF}}$ are simply the eigenvectors of \tilde{J} so that, in this case, β plays no role. Under the sparse regime, where $c = O_n(1)$, using the eigenvector associated to the smallest (resp., largest) eigenvalues of $H_{\beta, \tilde{J}}^{\text{MF}}$ (resp., \tilde{J}) as an estimator for σ does not allow to make non-trivial reconstruction as soon as theoretically possible, i.e. whenever $\beta_N > \beta_{\text{SG}} > \beta_F$: in this case indeed, the asymptotic spectrum of \tilde{J} is unbounded and no isolated eigenvalue of $H_{\beta, \tilde{J}}^{\text{MF}}$ is to be found. This explains the poor performance depicted in figure 6 for small average degrees c . On the opposite, as already observed in section 3.2, for sufficiently large degrees c , the naïve mean field approximation essentially yields the same result as the Bethe approximation.

4.3.3. The ‘spin glass Bethe-Hessian’. We conclude this section by presenting an alternative use of the Bethe-Hessian matrix, inspired by the work of [20], that we name here the *spin glass Bethe-Hessian*. Algorithm 2 represents an optimal relaxation of the Bayes optimal solution, capable of performing better than random inference as soon as theoretically possible. The parametrization $\beta = \beta_N$ is not the only possible choice of β able to reach this threshold. It was indeed shown, under different settings, in [20, 27, 46, 47] that choosing the temperature $\beta = \beta_{\text{SG}}$ allows one also to achieve non-trivial clustering as soon as theoretically possible.

The value β_{SG} , unlike β_{N} , can be easily estimated from the matrix \tilde{J} solving $c\mathbb{E}[\text{th}^2(\beta_{\text{SG}}\tilde{J}_{ij}\sigma_i\sigma_j)] = c\mathbb{E}[\text{th}^2(\beta_{\text{SG}}\tilde{J}_{ij})] = 1$. However, it was proved in [34] that for community detection in realistic *heterogeneous* (thus not Erdős–Rényi-like) graphs, this may be a quite suboptimal choice in terms of the raw (say, overlap) classification performance. The main difference between the *spin glass Bethe-Hessian* and the *Nishimori Bethe-Hessian* is thus observed when the underlying graph is not of an Erdős–Rényi type. This can be understood by a closer inspection of equation (28), which shows that the vector σ is an approximate eigenvector of $H_{\beta_{\text{N}},\tilde{J}}$ for *any* underlying degree distribution of the graph. This would not be true in general for any other value of $\beta \neq \beta_{\text{N}}$, hence in particular not for β_{SG} .

As a visual confirmation, figure 7 displays the overlap performance and the histograms of the entries of the informative eigenvector of $H_{\beta_{\text{N}},\tilde{J}}$ versus $H_{\beta_{\text{SG}},\tilde{J}}$ for a matrix \tilde{J} generated according to equation (26), considering on the top row graphs with an underlying power-law degree distribution (this thus goes beyond the assumption of the present article, yet is typical of real-world graph models [48]) and on the bottom row Erdős–Rényi graphs. The loss in precision of the *spin glass Bethe-Hessian* is best understood by comparing the two histograms which evidence that, unlike $H_{\beta_{\text{N}},\tilde{J}}$, the underlying node classes seen by $H_{\beta_{\text{SG}},\tilde{J}}$ is much spoiled by the heterogeneous degree distribution. This is also observed to some extent on Erdős–Rényi graphs, but here the performance achieved by $H_{\beta_{\text{SG}},\tilde{J}}$ is essentially the same as the one obtained with $H_{\beta_{\text{N}},\tilde{J}}$.

The use of $H_{\beta_{\text{N}},\tilde{J}}$ should thus be privileged when the input weighted graph \mathcal{G} may be far from an Erdős–Rényi random graph generation, such as in the case of a real-world weighted social graph. Besides, one can envision to extend algorithm 2 beyond two-class node clustering, as proposed in [43], where the authors show that the proper parametrization of the Bethe-Hessian matrix (specifically using multiple rather than a single value for β) brings a decisive advantage on real datasets.

On the opposite, if the input graph is of the Erdős–Rényi type, the performances of both algorithms are observed to be similar, with a slight computational as numerical stability advantage for $H_{\beta_{\text{SG}},\tilde{J}}$. We nonetheless underline that the estimation of β_{N} may be of independent interest: if one uses the solution of spectral clustering as the initialization to an algorithm seeking the actual Bayes optimal solution, then the initialization provided by $H_{\beta_{\text{SG}},\tilde{J}}$ would likely be of good quality, although β_{N} would still remain unknown.

4.4. Application to real data classification

We complete the article by a robustness test of our proposed algorithm under a real-world machine learning classification problem. Specifically, we consider a *sparse (and thus cost-efficient) version* of the problem of *correlation clustering* such as met in image classification and show how algorithm 2 can be adopted to accomplish this task with higher performance than with competing spectral methods of the literature.

Let $\{\mathbf{z}_i\}_{i=1,\dots,n}$ be an n -vector dataset with $\mathbf{z}_i \in \mathbb{R}^p$. These vectors represent discriminating *features* of some two-class data (say images) to be clustered in a fully unsupervised manner. In typical modern machine learning, p is of the order of a few

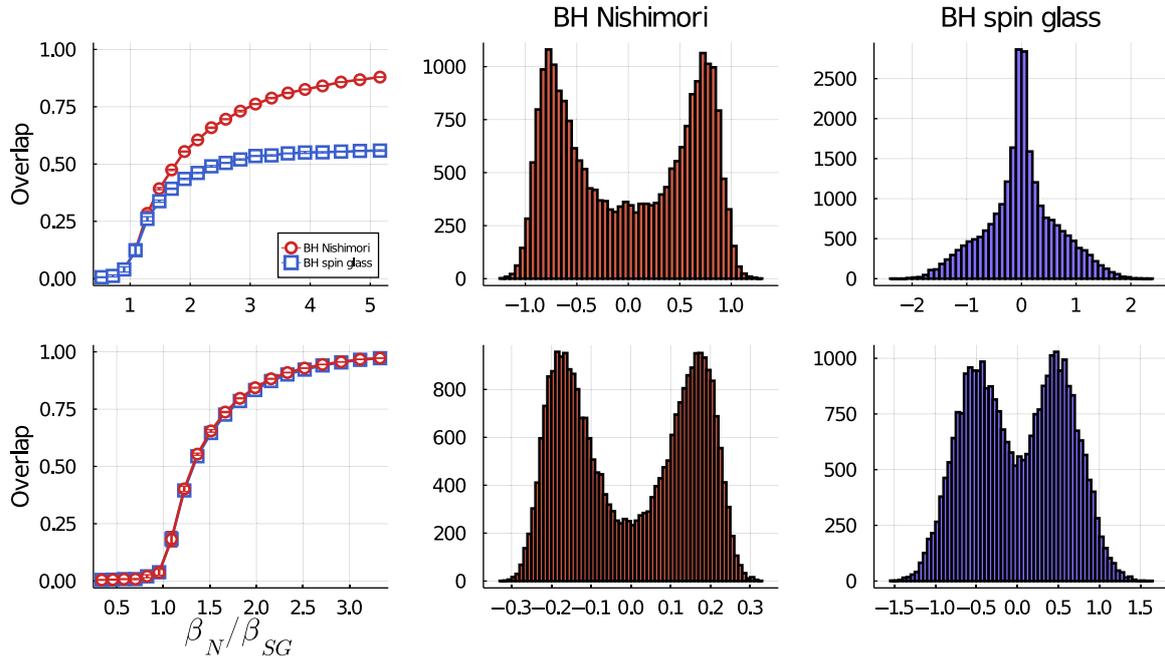


Figure 7. (First row) Random graphs with an underlying power law degree distribution. (Second row) Erdős Rényi random graphs. (First column) Overlap performance obtained exploiting the eigenvector associated to the smallest eigenvalue of $H_{\beta_N, \tilde{J}}$ (red circles) $H_{\beta_{SG}, \tilde{J}}$ (blue squares). The entries of \tilde{J} are distributed according to a Gaussian measure as in equation (26). Averages are taken over ten realizations. (Second column) Histogram of the entries of the informative eigenvector of $H_{\beta_N, \tilde{J}}$ for $\beta_N/\beta_{SG} \approx 3.6$ in the first plot. (Third column) Histogram of the entries of the informative eigenvector of $H_{\beta_{SG}, \tilde{J}}$ for the same configuration as the second plot. For all plots, the graphs have $n = 30\,000$ nodes and expected average degree $c = 10$.

thousands for images and a few hundreds for natural language text representations, and it is not rare to try and classify up to millions of data vectors \mathbf{z}_i .

The most elementary unsupervised machine learning classification approach consists in running the popular *k-means* algorithm in the ambient p -dimensional feature space. *K-means* is however known to fail for large p [49] and is ruled out as soon as p exceeds the order of a few tens. A classical workaround is to *embed* the feature vectors \mathbf{z}_i in a lower dimensional space on which to run *k-means* clustering. The most popular embedding exploits a *spectral* approach: one starts by defining a kernel matrix $K(\{\mathbf{z}\}) \in \mathbb{R}^{n \times n}$, the entry $K_{ij}(\{\mathbf{z}\})$ of which evaluates some *affinity metric* between \mathbf{z}_i and \mathbf{z}_j ; running a principal component analysis on $K(\{\mathbf{z}\})$, one then extracts a collection of eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ for some ℓ of the order of the presumed number of classes; the rows $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n \in \mathbb{R}^\ell$ of the resulting ‘tall’ matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_\ell] \in \mathbb{R}^{n \times \ell}$ form the embedding of the original features from \mathbb{R}^p into \mathbb{R}^ℓ over which *k-means* clustering is finally run. A popular affinity function is merely the correlation $K_{ij}(\{\mathbf{z}\}) = \mathbf{z}_i^T \mathbf{z}_j$, which we shall consider here⁹.

⁹Other choices exist, such as the more popular *heat* kernel $K_{ij}(\{\mathbf{z}\}) = \exp(-\|\mathbf{z}_i - \mathbf{z}_j\|^2/2\nu^2)$ for some $\nu > 0$.

For large dimensional datasets though (i.e. for p, n beyond a few thousands), the $O(pn^2)$ cost of building $K(\{\mathbf{z}\})$ added to the (at least) $O(n^2)$ cost of the principal component analysis step makes spectral clustering hardly achievable on a modern home computer. To drastically decrease the computational complexity one may proceed to a two-level sparsification as recently proposed in [50, 51]: by randomly discarding elements of the p -dimensional features \mathbf{z}_i and by randomly dropping a number of evaluations of the correlations $\mathbf{z}_i^T \mathbf{z}_j$. This operation of course impedes the clustering performance, but, as surprisingly proved in [50, 51] under a ‘still rather dense graph’ regime, the performance loss is negligible for a wide range of sparsity levels. To this end, let $S \in \{0, 1\}^{n \times p}$ and $M \in \{0, 1\}^{n \times n}$ (symmetric) be Bernoulli masks with parameters $\sqrt{\kappa/p}$ and c/n ,¹⁰ respectively. The resulting sparsified kernel matrix then becomes

$$\tilde{J} = K(\{\tilde{\mathbf{x}}\}) \circ M, \quad \text{where } \tilde{x}_{i,l} = x_i S_{i,l}. \quad (33)$$

i.e. each entry of each of the feature vectors \mathbf{z}_i is kept only with probability $\sqrt{\kappa/p}$, while each measurement $K_{ij}(\{\tilde{\mathbf{z}}\})$ is only performed with probability c/n . The computational complexity to build \tilde{J} is thus scaled down to $O(\kappa cn)$, i.e. to linear time complexity with respect to the size of the original dataset (which is the best one can hope for without completely dropping part some of the data \mathbf{z}_i).

As a major consequence of the sparsification procedure, the non-zero entries of \tilde{J} can be considered asymptotically independent due to the asymptotic absence of short loops in the underlying sparse Erdős–Rényi graph. As a result, equation (26) provides a good approximation for the generative model of \tilde{J} and for a two-class correlation clustering problem, algorithm 2 can be efficiently used on the matrix \tilde{J} .

We thus practically tested algorithm 2 against the naïve mean field approach which in this setting happens to coincide with the algorithm proposed in [50] when applied to the $\tilde{\mathbf{z}}_i$ vectors, and against the weighted Laplacian matrix approach. As a telling modern data classification context, we chose to cluster two classes of high-resolution extremely realistic images randomly produced by *generative adversarial networks* (the now quite popular GANs) [16]; the interest of using GAN images rather than real images lies in that GAN images can be produced ‘on-the-fly’ and in arbitrary numbers.

Specifically, we considered $n = 40\,000$ images divided into two groups of equal size, representing collie dogs and tabby cats. A representative example of the input images generated by the GAN is given in figure 8. For each of these images we extracted discriminating features using an off-the-shelf convolutional neural network (VGG) which produces $p = 512$ -dimensional feature vectors \mathbf{z}_i .¹¹ We then measured the overlap performance as a function of the average node degree c of the ensuing graph and for different values of κ . The results are reported in figure 8 which strikingly evidences that algorithm 2 can achieve almost perfect reconstruction already for $c = 5$ when the feature vectors \mathbf{z}_i are not sparsified ($\kappa = p$): so, in clearer terms, out of the $40k \times 40k = 1.6 \cdot 10^9$ correlations needed to evaluate the full $K(\{\mathbf{z}\})$ matrix, only $\approx 6 \times 40k = 2.4 \cdot 10^5$ is enough to achieve almost optimal performance, thus corresponding to a striking 10^4 -fold gain in complexity for a rather marginal performance loss!

¹⁰ The choice of c is not a coincidence: M will enforce an average node degree of c to the resulting graph.

¹¹ The $p = 512$ figure is on the low-hand of typical image vector representations: this number today may rise to $4k$ or even to $20k$ when much more than two classes of images are to be classified.

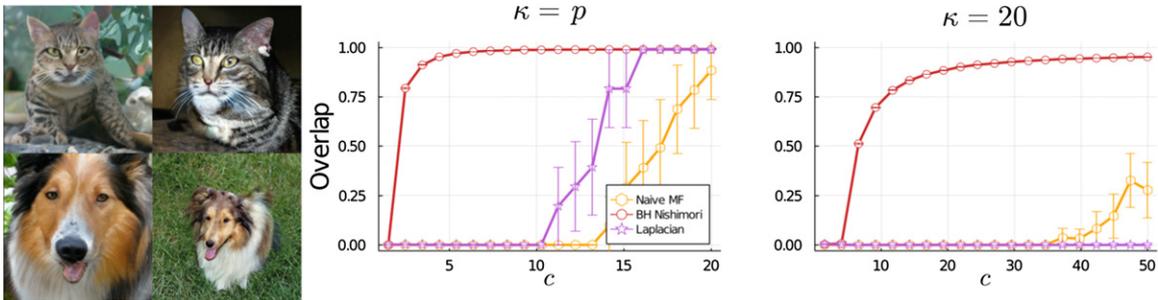


Figure 8. (Left) Example of random generations of GAN images representing collie dogs and tabby cats used for the experiment. (Middle and right) Overlap classification performance of 40 000 GAN images, as a function of the expected average underlying graph degree c . Here, we consider $K_{ij}(\{\mathbf{x}\}) = \frac{1}{p} \mathbf{x}_i^T \mathbf{x}_j$ and we take either $\kappa = p$: all features of the images are kept, or $\kappa = 20$: on average, only $\sqrt{\kappa/p}$ features (out of the original $p = 512$) are used. Simulation performances are themselves averaged over ten realizations.

Figure 8 also reports that the performance of the naïve mean-field and weighted Laplacian matrix approaches, currently the legacy methods in the literature, severely suffer in the low- c end. These observations perfectly adhere with the conclusions drawn so far in the article and thus turns our up-to-here formal *Nishimori-optimized* algorithm into a concrete powerful method for cost-efficient classification of large dimensional datasets.

5. Conclusion

The central contribution of the article is of a theoretical nature and aims at introducing an elegant explicit relation between the Bethe-Hessian matrix and the Nishimori temperature. Yet, beyond this statistical physics endeavour, which will surely find further independent theoretical interests, the result finds fundamental direct applications to Bayesian statistical inference; this is strikingly evidenced by the image clustering application devised in section 4.4. Specifically, one may anticipate an important impact in more involved applications than those considered in this article, such as in *restricted Boltzmann machines* (RBM) whose goal is to learn a generative model from a set of examples [52]: the Bethe approximation has recently been adopted to study the RBM from a Bayesian perspective [53] so that one may envision that the explicit relation between the Bethe free energy and the Bayes optimal (Nishimori) condition presented in this article would lead to a better understanding and improvement of state-of-the-art algorithms. Similarly, the Bethe and TAP approximations have recently been exploited to devise efficient spectral algorithms for phase retrieval, based on statistical physics intuitions similar to the ones detailed in this article [35, 54, 55]. The extension of our results to this more involved setting is a promising line of exploration.

On the side of complexity reduction, exploiting high levels of sparsification of data measurements, we showed that our proposed algorithm is capable of accomplishing high quality *unsupervised* classification on very large datasets. This result is all the more fundamental that future machine learning data treatment will call for increasingly larger datasets which cannot be possibly manually labeled and for which unsupervised (or possibly semi-supervised) approaches must be adopted¹². As a downside though, the generative model we considered for the data affinity (kernel) matrix takes the strong assumption that its entries are drawn from the same probability distribution and only the average (and not the variance, or the distribution itself) embeds information on the node labels. This setting might be too simplistic on generic real data that would require more realistic probability distributions for the generative model of the kernel matrix, considering, for instance, asymmetrical [20], multi-cluster [27] or multi-dimensional distributions.

Possibly most importantly, we worked here under the assumptions that the edges maintained in the sparsified graph are drawn independently at random. When dealing with actual kernel matrices, this cost-efficient measure is quite suboptimal: in [56], a more efficient sparsification procedure is used which maintains the entries of $K(\{\mathbf{z}\})$ of largest amplitude. In [56], this comes at the cost of computing all the entries of $K(\{\mathbf{z}\})$ but, surely, a more efficient nearest neighbors-type procedure could be implemented as a good performance-complexity compromise [57]. Yet, in this setting, although stronger sparsity levels can surely be achieved for the same performance, the key independence property of the entries of $K(\{\mathbf{z}\})$ which we exploited here can no longer be assumed, so that one needs to carefully handle the hard problem of dependencies. There lies the main objectives of our follow-up investigations.

Acknowledgments

R C's work is supported by the MIAI LargeDATA Chair at University Grenoble-Alpes and the GIPSA-HUAWEI Labs project Lardist. N T's work is partly supported by the French National Research Agency in the framework of the 'Investissements d'avenir' program (ANR-15-IDEX-02) and the LabEx PERSYVAL (ANR-11-LABX-0025-01). The authors thank Mohamed El Amine Seddik for sharing the codes to produce the experiments on GAN images.

Appendix A. An explicit expression for the matrix $F(\mathbf{g})$

We here provide one of the possible explicit expressions that the matrix $F(\mathbf{g})$ can have. In particular, this is the expression used in our simulations. Let us recall the definition of the matrix $F(\mathbf{g})$.

Let $\mathbf{g} \in \mathbb{R}^{2|\mathcal{E}|}$ be an eigenvector of the matrix B with weight vector $\boldsymbol{\omega} \in \mathbb{R}^{2|\mathcal{E}|}$. Let λ be the eigenvalue corresponding to \mathbf{g} , with $|\lambda| \geq 1$. The matrix $F(\mathbf{g})$ is *any* matrix

¹² A configuration which, in passing, even modern so-called *deep* neural networks struggle to correctly handle.

satisfying the relation

$$[F(\mathbf{g})\boldsymbol{\psi}(\mathbf{g})]_i = \sum_{j \in \partial i} \omega_{ij}^3 g_{ij}, \tag{34}$$

where we recall that

$$\psi_i(\mathbf{g}) = \sum_{j \in \partial i} \omega_{ij} g_{ij}.$$

A possible definition of the matrix $F(\mathbf{g})$ is to consider a diagonal matrix, satisfying

$$F_{ij}(\mathbf{g}) = \delta_{ij} \frac{\sum_{j \in \partial i} \omega_{ij}^3 g_{ij}}{\sum_{j \in \partial i} \omega_{ij} g_{ij}}.$$

This matrix depends however explicitly on \mathbf{g} . We here describe an alternative expression in which the dependence on \mathbf{g} is manifested only through λ . More explicitly, the following relation holds

$$\lambda g_{ij} = (B\mathbf{g})_{ij} = \psi_j(\mathbf{g}) - \omega_{ij} g_{ji}.$$

Considering the same equation for g_{ji} , we can easily write the following system

$$\begin{pmatrix} \lambda & \omega_{ij} \\ \omega_{ij} & \lambda \end{pmatrix} \begin{pmatrix} g_{ij} \\ g_{ji} \end{pmatrix} = \begin{pmatrix} \psi_j(\mathbf{g}) \\ \psi_i(\mathbf{g}) \end{pmatrix}.$$

For $|\lambda| \geq 1$ and $|\omega_{ij}| < 1$, the matrix on the left hand-side can be inverted, leading to the following relation

$$g_{ij} = \frac{\lambda \psi_j - \omega_{ij} \psi_i}{\lambda^2 - \omega_{ij}^2}. \tag{35}$$

Plugging equation (35) into equation (34), the following expression of $F(\mathbf{g}) \equiv F(\lambda)$ can be obtained:

$$F_{ij}(\lambda) = -\delta_{ij} \sum_{k \in \partial i} \frac{\omega_{ik}^4}{\lambda^2 - \omega_{ik}^2} + \frac{\lambda \omega_{ij}^3}{\lambda^2 - \omega_{ij}^2}.$$

This expression of the matrix $F(\mathbf{g}) \equiv F(\lambda)$ is the one considered in our simulations.

Appendix B. Algorithm implementation

In this appendix we discuss more extensively some details concerning a practical and efficient implementation of algorithm 2. For reference, our codes are available at github.com/lorenzodallamico/NishimoriBetheHessian. We now proceed to a detailed analysis of each step of algorithm 2.

The first step of algorithm 2 consists in the following operation:

$$\forall (ij) \in \mathcal{E} : \tilde{J}_{ij} = \tilde{J}_{ij} - \frac{1}{2|\mathcal{E}|} \mathbf{1}_n^T \tilde{J} \mathbf{1}_n.$$

The rationale of this operation is to consider an input matrix \tilde{J} as close as possible to a realization of the distribution of equation (26) that satisfy, for two classes of equal size¹³, the condition $\mathbb{E}[\tilde{J}_{ij}] = 0$. By shifting the empirical average of \tilde{J}_{ij} to zero for the input of algorithm 2, we are willing to reproduce this property. Note that *only the non-zero entries of \tilde{J} are shifted*, while for all the $(ij) \notin \mathcal{E}$ the $\tilde{J}_{ij} = 0$.

Once a proper input matrix \tilde{J} is obtained, the value of β_{SG} and then the smallest eigenvalue of $H_{\beta_{\text{SG}}, \tilde{J}}$ are computed: if the latter is positive, one cannot proceed any further to the computation of $\hat{\beta}_N$ and the algorithm is stopped. In the spirit of *correlation clustering*, discussed in section 4, the condition $\gamma_{\min}(H_{\beta_{\text{SG}}, \tilde{J}}) < 0$ imposes the minimal average degree to perform non-trivial reconstruction¹⁴.

At this point, we get to the core of algorithm 2 that consists in the computation of $\hat{\beta}_N$. The first thing to do is to determine if \mathcal{G} is a signed graph (with only $\pm J_0$ entries). If this is the case, the signed representation of $H_{\beta, \tilde{J}}$ introduced in equation (30) should be adopted. We consider first this easier case.

For notation convenience, we introduce $r = [\text{th}(\beta J_0)]^{-1}$ ($r \geq 1$) and define $H_{r, \tilde{J}} = (r^2 - 1)I_n + D - r\tilde{J}$. Furthermore, let \mathbf{x}_r be the eigenvector of $H_{r, \tilde{J}}$ associated to its smallest eigenvalues. We look for r so that

$$\gamma_{\min}(H_{r, \tilde{J}}) = 0.$$

In order to do so, consider $r_t > \hat{r}_N = [\text{th}(\hat{\beta}_N J_0)]^{-1}$. The following relation is true for any r ,

$$\gamma_{\min}(H_{r, \tilde{J}}) \leq \mathbf{x}_{r_t}^T H_{r, \tilde{J}} \mathbf{x}_{r_t} = (r^2 - 1) + d_{r_t} - r \mathbf{j}_{r_t} := f_{r_t}(r), \tag{36}$$

where $d_{r_t} = \mathbf{x}_{r_t}^T D \mathbf{x}_{r_t}$ and $\mathbf{j}_{r_t} = \mathbf{x}_{r_t}^T \tilde{J} \mathbf{x}_{r_t}$. Defining r_{t+1} as the solution to $f_{r_t}(r_{t+1}) = 0$, one immediately obtains from equation (36) that $\gamma_{\min}(H_{r_{t+1}, \tilde{J}}) < 0$. One can show (appendix F in [47]) that $|r_{t+1} - \hat{r}_N| < |r_t - \hat{r}_N|$, hence, that at each iteration the value of r_t approaches \hat{r}_N . In practice, convergence is typically achieved in less than ten iterations. A good initialization is $r_0 = [\text{th}(\beta_{\text{SG}} J_0)]^{-1} > \hat{r}_N$ (recall that $\beta_N > \beta_{\text{SG}}$), ensuring the algorithm convergence.

We now consider graphs with non-binary weights that introduce additional complications. The entries of $H_{\beta, \tilde{J}}$ grow exponentially, with β , making the eigenvalue computation potentially unstable. In order to work with a matrix with entries of order 1, we introduce

¹³Note that the inference problem of equation (27) does not make any assumption on the respective sizes of the classes that can therefore be arbitrary. In the case of asymmetric classes, however, the term $\mathbb{E}[J_{ij}] \neq 0$ depends on the sizes of the two classes. In order to do the proper shift, one would then need additional information on the class sizes.

¹⁴The detectability condition we recall to imposed by $\beta_N > \beta_{\text{SG}}$. While β_N is independent of the average degree, β_{SG} is a decreasing function of the average degree, as it can be easily obtained from its definition in equation (8).

the following *weighted regularized Laplacian* [58]:

$$L_{\beta, \tilde{J}} = I_n - \Lambda_{\beta, \tilde{J}}^{-1/2} \tilde{W}_{\beta, \tilde{J}} \Lambda_{\beta, \tilde{J}}^{-1/2}, \quad (37)$$

where

$$\left(\tilde{W}_{\beta, \tilde{J}}\right)_{ij} = \frac{\text{th}(\beta \tilde{J}_{ij})}{1 - \text{th}^2(\beta \tilde{J}_{ij})}; \quad \left(\Lambda_{\beta, \tilde{J}}\right)_{ij} = \delta_{ij} \sum_{k \in \partial i} \frac{\text{th}^2(\beta \tilde{J}_{ik})}{1 - \text{th}^2(\beta \tilde{J}_{ik})}$$

It is straightforward to see that if $H_{\beta_N, \tilde{J}} \mathbf{x} = 0$, then $L_{\beta_N, \tilde{J}} \mathbf{v} = 0$, where $\mathbf{v} = \Lambda_{\beta, \tilde{J}}^{-1/2} \mathbf{x}$. The matrix $L_{\beta, \tilde{J}}$ hence allows to compute $\hat{\beta}_N$ and \mathbf{x} in a more efficient way, since it is more suited to eigenvalue computations. We can define in this case $f_{\beta_t}(\beta) = \mathbf{v}_t^T L_{\beta, \tilde{J}} \mathbf{v}_t$ and update β_{t+1} as the solution to $f_{\beta_t}(\beta_{t+1}) = 0$.

In any case, for very large values of β_N (hence for very easy clustering problems) numerical instabilities may occur. In order to avoid this problem, we allow a ‘maximal value’ β_{th} for $\hat{\beta}_N$ beyond which the algorithm is stopped. The main reason that allows us to do so is that if $\beta_N > \beta_{\text{th}}$ we are practically in an easy detection regime, for which the knowledge of the exact value of β_N is less relevant and can be otherwise achieved first estimating the labels $\hat{\sigma}$ (the estimation of which will be very accurate) and then solving $\mathbb{E}[\text{th}(\beta_N \tilde{J}_{ij} \sigma_i \sigma_j)] = \mathbb{E}[\text{th}^2(\beta_N \tilde{J}_{ij})]$. We empirically observed that a good stopping criterion is obtained imposing $\beta_{\text{th}} \sim \sqrt{c} \beta_{\text{SG}}$.

References

- [1] Binder K and Young A P 1986 Spin glasses: experimental facts, theoretical concepts, and open questions *Rev. Mod. Phys.* **58** 801
- [2] Jordan M I 1998 *Learning in Graphical Models* vol 89 (Berlin: Springer)
- [3] Wainwright M J and Jordan M I 2008 *Graphical Models, Exponential Families, and Variational Inference* (Hanover, MA: Now Publishers Inc.) (<https://doi.org/10.1561/22000000001>)
- [4] Opper M and Saad D 2001 *Advanced Mean Field Methods: Theory and Practice* (Cambridge, MA: MIT Press)
- [5] Nishimori H 2001 *Statistical Physics of Spin Glasses and Information Processing: An Introduction* vol 111 (Oxford: Clarendon)
- [6] Zdeborová L and Krzakala F 2016 Statistical physics of inference: thresholds and algorithms *Adv. Phys.* **65** 453–552
- [7] Nishimori H 1981 Internal energy, specific heat and correlation function of the bond-random Ising model *Prog. Theor. Phys.* **66** 1169–81
- [8] Nishimori H and Sherrington D 2001 Absence of replica symmetry breaking in a region of the phase diagram of the Ising spin glass *AIP Conf. Proc.* **553** 67–72
- [9] Georges A, Hansel D, Le Doussal P and Bouchaud J-P 1985 Exact properties of spin glasses: II. Nishimori’s line: new results and physical implications *J. Phys. France* **46** 1827–36
- [10] Gruzberg I A, Read N and Ludwig A W W 2001 Random-bond Ising model in two dimensions: the Nishimori line and supersymmetry *Phys. Rev. B* **63** 104422
- [11] Parisen Toldin F, Pelissetto A and Vicari E 2009 Strong-disorder paramagnetic–ferromagnetic fixed point in the square-lattice $\pm J$ Ising model *J. Stat. Phys.* **135** 1039–61
- [12] Iba Y 1999 The Nishimori line and Bayesian statistics *J. Phys. A: Math. Gen.* **32** 3875
- [13] Bansal N, Blum A and Chawla S 2004 Correlation clustering *Mach. Learn.* **56** 89–113
- [14] Langone R, Mall R, Alzate C and Suykens J A K 2016 Kernel spectral clustering and applications *Unsupervised Learning Algorithms* (Berlin: Springer) pp 135–61
- [15] Watanabe Y and Fukumizu K 2009 Graph zeta function in the Bethe free energy and loopy belief propagation *Advances in Neural Information Processing Systems* vol 22 pp 2017–25

- [16] Brock A, Donahue J and Simonyan K 2018 Large scale GAN training for high fidelity natural image synthesis (arXiv:1809.11096)
- [17] Edwards S F and Anderson P W 1975 Theory of spin glasses *J. Phys. F: Met. Phys.* **5** 965
- [18] Mezard M and Montanari A 2009 *Information, Physics, and Computation* (Oxford: Oxford University Press)
- [19] Thouless D J 1986 Spin-glass on a Bethe lattice *Phys. Rev. Lett.* **56** 1082
- [20] Saade A, Lelarge M, Krzakala F and Zdeborová L 2016 Clustering from sparse pairwise measurements *IEEE Int. Symp. Information Theory (ISIT)* (Piscataway, NJ: IEEE) pp 780–4
- [21] Krzakala F, Moore C, Mossel E, Neeman J, Sly A, Zdeborová L and Zhang P 2013 Spectral redemption in clustering sparse networks *Proc. Natl Acad. Sci.* **110** 20935–40
- [22] Zhang P 2015 Nonbacktracking operator for the Ising model and its applications in systems with multiple states *Phys. Rev. E* **91** 042120
- [23] Aleja D, Criado R, García del Amo A J, Pérez Á and Romance M 2019 Non-backtracking pagerank: from the classic model to Hashimoto matrices *Chaos Solitons Fractals* **126** 283–91
- [24] Torres L, Suárez-Serrato P and Eliassi-Rad T 2019 Non-backtracking cycles: length spectrum theory and graph mining applications *Appl. Netw. Sci.* **4** 41
- [25] Torres L, Chan K S, Tong H and Eliassi-Rad T 2020 Node immunization with non-backtracking eigenvalues (arXiv:2002.12309)
- [26] Arrigo F, Higham D J and Noferini V 2020 Beyond non-backtracking: non-cycling network centrality measures *Proc. R. Soc. A* **476** 20190653
- [27] Shi C, Liu Y and Zhang P 2018 Weighted community detection and data clustering using message passing *J. Stat. Mech.* **033405**
- [28] Watanabe Y and Fukumizu K 2011 Loopy belief propagation, Bethe free energy and graph zeta function (arXiv:1103.0605)
- [29] Sato I, Mitsuhashi H and Morita H 2014 A matrix-weighted zeta function of a graph *Linear Multilinear Algebra* **62** 114–25
- [30] Gulikers L, Lelarge M and Massoulié L 2016 Non-backtracking spectrum of degree-corrected stochastic block models (arXiv:1609.02487)
- [31] Bordenave C, Lelarge M and Massoulié L 2015 Non-backtracking spectrum of random graphs: community detection and non-regular Ramanujan graphs *IEEE 56th Annual Symp. Foundations of Computer Science* (Piscataway, NJ: IEEE) pp 1347–57
- [32] Stephan L and Massoulié L 2020 Non-backtracking spectra of weighted inhomogeneous random graphs (arXiv:2004.07408)
- [33] Coste S and Zhu Y 2019 Eigenvalues of the non-backtracking operator detached from the bulk (arXiv:1907.05603)
- [34] Dall’Amico L, Couillet R and Tremblay N 2019 Revisiting the Bethe-Hessian: improved community detection in sparse heterogeneous graphs *Advances Neural Information Processing Systems* pp 4039–49
- [35] Maillard A, Krzakala F, Lu Y M and Zdeborová L 2020 Construction of optimal spectral methods in phase retrieval (arXiv:2012.04524)
- [36] Bollobás B and Béla B 2001 *Random Graphs* vol 73 (Cambridge: Cambridge University Press)
- [37] Horton M D, Stark H M and Terras A A 2006 What are zeta functions of graphs and what are they good for? *Contemp. Math.* **415** 173–89
- [38] Sylvester J R 2000 Determinants of block matrices *Math. Gaz.* **84** 460–7
- [39] Wigner E P 1958 On the distribution of the roots of certain symmetric matrices *Ann. Math.* **67** 325–7
- [40] Bauer F L and Fike C T 1960 Norms and exclusion theorems *Numer. Math.* **2** 137–41
- [41] Mézard M, Parisi G and Virasoro M A 1987 *Spin glass theory and beyond: An Introduction to the Replica Method and Its Applications* vol 9 (Singapore: World Scientific)
- [42] Saad Y 1992 *Numerical Methods for Large Eigenvalue Problems* (Manchester: Manchester University Press)
- [43] Dall’Amico L, Couillet R and Tremblay N 2020 A unified framework for spectral clustering in sparse graphs (arXiv:2003.09198)
- [44] Von Luxburg U 2007 A tutorial on spectral clustering *Stat. Comput.* **17** 395–416
- [45] Kunegis J, Schmidt S, Lommatzsch A, Lerner J, De Luca E W and Albayrak S 2010 Spectral analysis of signed graphs for clustering, prediction and visualization *Proc. 2010 SIAM Int. Conf. Data Mining* (Philadelphia, PA: SIAM) pp 559–70
- [46] Saade A, Krzakala F and Zdeborová L 2014 Spectral clustering of graphs with the Bethe hessian (arXiv:1406.1880)
- [47] Dall’Amico L, Couillet R and Tremblay N 2020 Community detection in sparse time-evolving graphs with a dynamical Bethe-Hessian *Advances in Neural Information Processing Systems* vol 33 (Red Hook, NY: Curran Associates, Inc.) pp 7486–97

- [48] Barabási A-L and Albert R 1999 Emergence of scaling in random networks *Science* **286** 509–12
- [49] Krieger H-P, Kröger P and Zimek A 2009 Clustering high-dimensional data *ACM Trans. Knowl. Discov. Data* **3** 1–58
- [50] Zarrouk T, Couillet R, Chatelain F and Le Bihan N 2020 Performance-complexity trade-off in large dimensional statistics *IEEE 30th Int. Workshop on Machine Learning for Signal Processing (MLSP)* (Piscataway, NJ: IEEE) pp 1–6
- [51] Couillet R, Chatelain F and Le Bihan N 2021 Two-way kernel matrix puncturing: towards resource-efficient PCA and spectral clustering *38th Int. Conf. Machine Learning* pp 2156–65
- [52] Ackley D H, Hinton G E and Sejnowski T J 1985 A learning algorithm for Boltzmann machines *Cogn. Sci.* **9** 147–69
- [53] Huang H and Toyozumi T 2016 Unsupervised feature learning from finite data by message passing: discontinuous versus continuous phase transition *Phys. Rev. E* **94** 062310
- [54] Luo W, Alghamdi W and Lu Y M 2019 Optimal spectral initialization for signal recovery with applications to phase retrieval *IEEE Trans. Signal Process.* **67** 2347–56
- [55] Ma J, Dudeja R, Xu J, Maleki A and Wang X 2021 Spectral method for phase retrieval: an expectation propagation perspective *IEEE Trans. Inf. Theory* **67** 1332–55
- [56] Liao Z, Couillet R and Mahoney M W 2020 Sparse quantized spectral clustering (arXiv:2010.01376)
- [57] Muja M and Lowe D G 2009 Fast approximate nearest neighbors with automatic algorithm configuration *VISAPP* vol 2
- [58] Dall’Amico L, Couillet R and Tremblay N 2020 Optimal Laplacian regularization for sparse spectral community detection *ICASSP 2020–2020 IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP)* (Piscataway, NJ: IEEE) pp 3237–41